

1 **Somatic mutations and clonal dynamics in healthy and cirrhotic human**
2 **liver**

3

4 **Authors:**

5 Simon F Brunner (1); Nicola D Roberts (1); Luke A Wylie (1); Luiza Moore (1);
6 Sarah J Aitken (2,3); Susan E Davies (3); Mathijs A Sanders (1,4); Pete Ellis (1);
7 Chris Alder (1); Yvette Hooks (1); Federico Abascal (1); Michael R Stratton (1);
8 Inigo Martincorena (1); Matthew Hoare (2,5) *; Peter J Campbell (1,6) *.

9

10 * Co-corresponding authors

11

12 **Institutes:**

13 (1) Cancer Genome Project, Wellcome Trust Sanger Institute, Hinxton, CB10
14 1SA, UK

15 (2) CRUK Cambridge Institute, Robinson Way, Cambridge, CB2 0RE, UK

16 (3) Department of Pathology, University of Cambridge, Addenbrooke's
17 Hospital, Cambridge, CB2 0QQ, UK

18 (4) Department of Hematology, Erasmus University Medical Center,
19 Rotterdam, The Netherlands

20 (5) Department of Medicine, University of Cambridge, Addenbrooke's
21 Hospital, Cambridge, CB2 0QQ, UK

22 (6) Department of Haematology and Stem Cell Institute, University of
23 Cambridge, Hills Rd, Cambridge CB2 0XY, UK

24

25 **Address for correspondence:**

26 Dr Peter J. Campbell, Cancer Genome Project, Wellcome Trust Sanger Institute,
27 Hinxton CB10 1SA, United Kingdom. Telephone: +44 (0) 1223 834244. e-mail:
28 pc8@sanger.ac.uk

29

30 Dr Matthew Hoare, CRUK Cambridge Institute, Robinson Way, Cambridge, CB2
31 0RE United Kingdom. e-mail: Matthew.Hoare@cruk.cam.ac.uk

32

33 **SUMMARY**

34 The commonest causes of chronic liver disease are excess alcohol intake, viral
35 hepatitis or non-alcoholic fatty liver disease, with the clinical spectrum ranging
36 in severity from hepatic inflammation through cirrhosis to liver failure or
37 hepatocellular carcinoma. The hepatocellular carcinoma genome exhibits diverse
38 mutational signatures, resulting in recurrent mutations across >20-30 cancer
39 genes¹⁻⁷. Stem cells from normal livers have low mutation burden and limited
40 diversity of signatures⁸, suggesting that the complexity of hepatocellular
41 carcinoma arises during progression to chronic liver disease and subsequent
42 malignant transformation. We sequenced whole genomes of 482
43 microdissections of 100-500 hepatocytes from 5 normal and 9 cirrhotic livers.
44 Compared to normal liver, cirrhotic liver had higher mutation burden. Although
45 rare in normal hepatocytes, structural variants, including chromothripsis, were
46 prominent in cirrhosis. Driver mutations, both point mutations and structural
47 variants, affected 1-5% clones. Clonal expansions millimetres in diameter
48 occurred in cirrhosis, sequestered by bands of fibrosis engirdling regenerative
49 nodules. Some mutational signatures were universal and equally active in both
50 non-malignant hepatocytes and HCC; some were substantially more active in
51 HCC than chronic liver disease; and others, arising from exogenous exposures,
52 were present in a subset of patients. Up to 10-fold within-patient variation in
53 activity of exogenous signatures existed between adjacent cirrhotic nodules,
54 arising from clone-specific and microenvironmental forces. Synchronous
55 hepatocellular carcinomas exhibited the same mutational signatures as
56 background cirrhotic liver, but with higher burden. Somatic mutations chronicle
57 the exposures, toxicity, regeneration and clonal structure of liver tissue as it
58 progresses from health to disease.

59

60 **MAIN TEXT**

61 Identifying somatic mutations in non-malignant tissue requires approaches to
62 overcome its polyclonality, such as single cell sequencing⁹, cultures of single
63 cells^{8,10} or microbiopsy sequencing¹¹. The latter relies on local cell division with
64 limited migration leading to a clonal patchwork, a known property of
65 hepatocytes¹². We generated whole genome sequences from 482 laser-capture
66 microdissections of 100-500 hepatocytes (**Extended Figure 1A**) across 14
67 patients: 5 normal controls; 4 with cirrhosis from alcohol-related liver disease
68 (ARLD) and 5 with cirrhosis from non-alcoholic fatty liver disease (NAFLD)
69 (**Supplementary Tables 1-2, Extended Figures 4-6**). Samples of normal liver
70 were acquired from hepatic resections of colorectal cancer metastases; samples
71 of cirrhotic liver from patients transplanted for synchronous but distant
72 hepatocellular carcinoma (HCC).

73

74 To evaluate sensitivity and specificity, we generated independent libraries and
75 sequencing data from different sections of the same biopsy, microdissecting the
76 same x,y-region from adjacent z-stacks, separated by ~20 μ m. Concordance was
77 high between variants called in adjacent sections, but not distant pairs,
78 suggesting that specificity of mutation calls was high (**Extended Figure 1B**), and
79 sensitivity across patients was 50-95%, dependent on coverage and clonality
80 (**Extended Figure 1C-F**). As a further check on specificity, deep targeted
81 sequencing of cancer genes in the same library as 96 whole-genome samples
82 confirmed 16 of 17 mutations originally called. In keeping with polyploidy as a
83 late differentiation stage in liver¹³, 20-25% of mature hepatocytes in
84 microdissected samples were multinuclear (**Extended Figure 1G**). We therefore
85 deployed copy number algorithms with expected ploidy of 4, and report
86 mutation burdens per diploid genome, rather than per cell.

87

88 We observed considerable heterogeneity in burden of somatic substitutions both
89 between and within patients (**Figure 1A; Supplementary Tables 3-4**). Using
90 mixed effects models, microdissections from cirrhotic livers had, on average,

91 1251 (CI_{95%} 233-2268; p=0.02) extra substitutions per diploid genome
92 compared to normal livers, independent of age. In accordance with published
93 values⁸, the estimated rate of mutation accumulation was 33/year/diploid
94 genome, albeit with wide confidence intervals (CI_{95%} -17-84; p=0.18) and
95 moderate variation between individuals (estimated between-individual SD,
96 13/year). Indels showed the same heterogeneity between and within individuals
97 as substitutions (**Figure 1B**).

98

99 Structural variants and copy number alterations occurred in moderate numbers
100 across all 9 patients with liver cirrhosis, despite being rare in normal liver
101 (**Figure 1C, Extended Figure 2, Supplementary Tables 3-4**). Occasional whole
102 chromosome or arm-level aneuploidy occurred, as well as focal events, including
103 deletions, tandem duplications and unbalanced translocations (**Extended Figure**
104 **2**). We found 5 separate clusters of SVs, across 3 patients, with patterns
105 indicative of chromothripsis¹⁴ (**Figures 1D-F, Extended Figure 2**).
106 Chromothripsis, in which multiple rearrangements occur in a single catastrophic
107 mitosis¹⁴, is a major mutational process in cancers, occurring in ~5% of HCCs¹⁵,
108 but is rare in normal somatic cells. To see 1-2% of clones in chronic liver disease
109 with chromothripsis suggests that sustained toxicity and regeneration
110 substantially increases mitotic stress in hepatocytes.

111

112 We screened for driver mutations among coding regions, 5'-UTRs, 3'-UTRs and
113 promoters (**Supplementary Tables 5-8**). No elements were significant after
114 genome-wide multiple hypothesis correction, so we focused on the 30 most
115 prevalent HCC genes¹⁻⁵. These carried 22 non-synonymous variants, seen in both
116 normal and cirrhotic samples, including inactivating mutations in the tumour
117 suppressor genes *ACVR2A*, *ARID2*, *ARID1A* and *TSC2* (**Extended Figure 3A**). With
118 hypothesis testing restricted to these 30 genes, *ALB* (q=0.001) and *ACVR2A*
119 (q=0.001) were significant. Recurrence in *ALB* (albumin) likely reflects a
120 mutational process in which indels preferentially occur in highly expressed
121 genes, as reported in HCCs^{5,16} (**Extended Figure 3B-C**). Assuming no negative
122 selection, we can use the ratio of non-synonymous to synonymous substitutions

123 for the 30 HCC genes to estimate the number of driver substitutions among
124 them¹⁷ – this gives a 95% confidence interval of 0.0–13.2 drivers in total across
125 482 microdissections (<3%). Among copy number aberrations of potential
126 significance^{1,2,18} (**Supplementary Table 9**), we found instances of chromosome
127 22 loss, 8q gain and 8p loss. Two focal deletions in different patients spanned
128 *ACVR2A* (**Extended Figure 2C,E**). We also found a reciprocal inversion that
129 deleted *CDKN2A* (**Extended Figure 2F**), the most common focal deletion in HCC,
130 and a deletion affecting *ARID5A*.

131

132 We reconstructed phylogenetic trees¹⁹, layering them onto the specimen's
133 histology. Samples from the healthy controls showed the highly polyclonal
134 nature of normal liver, with little genetic relatedness among even closely located
135 microdissections (**Figure 2A-D, Extended Figure 4**). Samples from patients
136 with chronic liver disease showed more complex clonal structure, from which
137 three general inferences can be drawn (**Figure 2E-P, Extended Figures 5-6**).
138 First, we found no sharing of mutations between adjacent liver nodules
139 separated by fibrotic bands. This suggests that the connective tissue laid down
140 during cycles of damage and regeneration sequesters clones from early stages of
141 the disease process. Second, some cirrhotic nodules were monoclonally derived
142 (**Figure 2J,N**, for example), while others were oligoclonal (**Figure 2F**), with
143 shared mutations often extending across microdissections millimetres apart.
144 Third, branching structures in phylogenies point to subclonal diversification
145 within nodules. Within such a clone, the proportion of shared, clonal mutations
146 on the trunk relative to those on the subclonal branches gives an estimate in
147 molecular time of when the most recent common ancestor of the clone emerged.
148 In some patients (for example, **Figure 2I-J**), the common ancestor of individual
149 nodules emerged relatively early in molecular time, while in others (**Figure 2M-**
150 **N**), the common ancestor appeared much more recently. Since the majority of
151 liver cells do not have driver mutations, the size and rapidity of clonal
152 expansions observed here evince the considerable in-built capacity of
153 hepatocytes to regenerate in response to liver damage.

154

155 A major debate in modelling cancer development is whether cancers need higher
156 mutation rates in order to acquire sufficient drivers. We compared mutation
157 burden in cirrhotic liver to synchronous, clonally unrelated HCCs from 7
158 patients. Synchronous HCCs carried, on average, 4600 more mutations than
159 matched cirrhotic liver (CI_{95%} 3600-5500; $p < 10^{-18}$ LME models; **Figure 3A**). This
160 argues that mutation rates increase during malignant transformation, either
161 through cancer-specific mutational processes or through greater activity in
162 cancers of widespread mutational processes.

163

164 To assess what mutational processes are active in cirrhosis, we extracted
165 mutational signatures across our 482 microdissections, the 7 synchronous HCCs
166 and 54 HCC genomes from TCGA¹, using two independent algorithms (**Figure**
167 **3B-E, Extended Figures 7-8**). Three major groups of mutational signatures
168 emerged: those ubiquitous and similarly active across cirrhosis and HCC; those
169 quiet in cirrhosis but universally more active in HCC; and those contributing to
170 some patients but not others, including signatures arising from exogenous
171 exposures.

172

173 In normal and cirrhotic liver, ubiquitous mutational signatures (5 and Sig.A)
174 were prevalent across clones, typically accounting for >75% of mutations in
175 combination. Signature 5 is widespread across cancers, including HCCs^{2,4,20}, and
176 accumulates linearly with age, suggesting it arises from endogenous mutational
177 processes. Sig.A is the dominant cause of mutations in normal blood stem
178 cells^{10,21} and leukaemias²¹, suggesting it too arises endogenously. In HCCs,
179 although Sig.A accounted for a lower proportion of mutations than in normal or
180 cirrhotic liver, the absolute numbers of mutations attributed to Sig.A were
181 comparable (Difference between cancer and non-cancer, 60 mutations; CI_{95%} -
182 80-200; $p = 0.4$; **Figure 3F, Supplementary Table 10**). This suggests that it is
183 active in hepatocytes throughout life, but is outstripped in HCC by mutational
184 processes emerging during malignant transformation.

185

186 A second group of mutational signatures comprises processes that are relatively
187 quiet in cirrhotic liver but universally more active in HCC (signatures 1, 12, 16,
188 40 and a novel signature, D; **Supplementary Table 10**). One of these, signature
189 16, consists of T>C mutations in ApT context and has a known transcriptional
190 strand bias, with both preferential repair of damaged adenines on transcribed
191 strands and increased damage on non-transcribed strands²². Although this
192 signature is more active in HCCs, we do see its characteristic transcriptional
193 strand bias in cirrhotic liver (**Extended Figure 9A**). Signature 1, caused by
194 spontaneous deamination of methylated cytosine to thymine, is also much more
195 active in HCC than non-malignant liver. The acceleration and universality of
196 these signatures in HCC suggests they reflect inbuilt DNA damage and repair
197 processes in hepatocytes that are unmasked during malignant transformation.

198

199 The third group of mutational processes represents signatures seen sporadically
200 across the cohort, many of which are due to exogenous exposures. One, signature
201 4, is found in lung cancers from smokers²⁰ and also HCCs, albeit with a less clear-
202 cut relationship to tobacco². Of our 14 patients, 4 had >10% of microdissections
203 with >5% of mutations attributed to signature 4, showing the expected
204 transcriptional strand bias on guanines (**Extended Figure 9B**). Not only did
205 signature 4 show considerable patient-to-patient heterogeneity, there was also
206 unexpectedly high clone-to-clone and nodule-to-nodule variability within
207 individual livers. In one patient, for example, about half the clones we sequenced
208 had 2000-4000 mutations, whereas the other half had 8000-12000, driven by
209 presence or absence of signature 4 (**Figure 4A**).

210

211 This within-patient regional variability extended to other exogenous exposures.
212 In one patient, 20-35% of mutations derived from signature 22 (**Figure 4B**;
213 **Extended Figure 9C**), characteristic of exposure to aristolochic acid²³. This
214 patient grew up in Poland, holidaying in Balkan states where aristolochic acid
215 exposure is pervasive²⁴. In a different patient, a subset of microdissections had
216 10-20% mutations attributable to signature 24 (**Figure 4C**), associated with
217 aflatoxin-B₁ exposure⁵. Biomarkers of exposure to aflatoxin-B₁, produced by

218 *Aspergillus* moulds contaminating crops, are prevalent in arable farmers²⁵, the
219 occupation of our patient. In both patients, these carcinogens showed striking
220 variability in mutational activity over short distances, generating few mutations
221 in some clones and hundreds to thousands in others – such striking regional
222 variation in activity of exogenous signatures is both unexpected and
223 unexplained.

224

225 In one patient, we found a large clone that carried >2000 mutations attributed to
226 signature 9 (**Figure 4D**), caused by off-target somatic hypermutation in B
227 lymphocytes²⁰. A clonotypic *IGH* rearrangement was evident, consistent with a
228 single B lymphocyte subclonally diversifying as it expanded in the liver
229 (**Extended Figure 10**). Signature 9 was only present on the ancestral trunk,
230 whereas signatures in the subclones, acquired in the liver, distributed similarly
231 to hepatocytes, suggesting the hepatic microenvironment shaped the on-going
232 mutational processes in the lymphocytes.

233

234 In conclusion, then, non-malignant liver has considerably lower proportions of
235 clones (<5%) with driver point mutations or structural variants than oesophagus
236 or skin^{11,26,27}, and those present were seen in both normal and cirrhotic liver.
237 They did not drive large clonal expansions, being restricted by fibrosis, and were
238 not shared with the distant synchronous HCCs, suggesting that the increased
239 cancer risk seen in chronic liver disease arises from a myriad of clones
240 competing independently to acquire sufficient driver mutations. *TERT* promoter
241 mutations are likely to be key events in this progression as they are seen in
242 dysplastic hepatic nodules^{18,28}, but we did not identify any in cirrhotic or normal
243 liver. The low proportion of clones with drivers observed here and in exome
244 studies performed elsewhere^{29,30} means that much larger sample sizes will be
245 needed to comprehensively map how driver mutations accumulate in the
246 progression from normal liver through regenerative and dysplastic nodules to
247 HCC.

248

249 These data reveal the genomic consequences of chronic liver disease – increased
250 mutation rates; complex structural variation including chromothripsis;
251 aneuploidies; low burden of mutations targeting known HCC genes. Genomically,
252 one middle-aged, healthy liver looks much like any other: a community of small,
253 tightly packed clones, each comprising a few hundred cells, containing ~1000-
254 1500 mutations, painted from a limited palette of signatures. Unhealthy livers
255 diverge from this norm: large dynasties of clones, sequestered by impassable
256 bands of fibrosis, their palette of signatures more variable, more vigorous, more
257 regionally variegated.

258

259 **ACKNOWLEDGEMENTS**

260 This work was supported by a the Wellcome Trust and a Cancer Research UK
261 Grand Challenge Award [C98/A24032]. PJC is a Wellcome Trust Senior Clinical
262 Fellow (WT088340MA). SFB was supported by the Swiss National Science
263 Foundation (P2SKP3-171753 and P400PB-180790). MAS is supported by a
264 Rubicon fellowship from NWO (019.153LW.038). The Cambridge Human
265 Research Tissue Bank is supported by the NIHR Cambridge Biomedical Research
266 Centre. MH is supported by a CRUK Clinician Scientist Fellowship
267 (C52489/A19924).

268

269 **DATA AVAILABILITY**

270 Whole genome sequencing data in the form of BAM files across samples reported
271 in this study have been deposited in the European Genome-Phenome Archive
272 (<https://www.ebi.ac.uk/ega/home>) with accession number
273 EGAD00001004578. Substitution and indel calls have been deposited on
274 Mendeley Data ('Somatic mutations and clonal dynamics in healthy and cirrhotic
275 human liver': <http://dx.doi.org/10.17632/ktx7jp8sch.1>).

276

277 **CODE AVAILABILITY**

278 Single-nucleotide substitutions were called using the CaVEMan (cancer variants
279 through expectation maximization) algorithm, version 1.11.2
280 (<https://github.com/cancerit/CaVEMan>). Small insertions and deletions were
281 called using the Pindel algorithm, version 2.2.2

282 (<https://github.com/genome/pindel>). Rearrangements were called using the
283 BRASS (breakpoint via assembly) algorithm version 5.4.1
284 (<https://github.com/cancerit/BRASS>). Miscellaneous scripts for downstream
285 analysis are available on Github (<https://github.com/sfbrunner/liver-pub-repo>).
286 Mutational signatures analysis performed using the HDP hierarchical Dirichlet
287 Process package version 0.1.5, available on Github
288 (<https://github.com/nicolaroberts/hdp>).

289

290 **AUTHOR CONTRIBUTIONS**

291 P.J.C., M.H. and S.F.B. designed the experiments. S.F.B. performed the laser-
292 capture microdissection, data curation and statistical analysis, with L.A.W.,
293 M.A.S., F.A. and I.M. providing assistance and advice. M.H., S.J.A. and S.E.D.
294 collated and analysed the clinical and histological data from the patients. N.D.R.
295 developed the hierarchical Dirichlet process for extracting mutational
296 signatures. L.M. and P.E. developed the laser-capture microdissection, DNA
297 extraction and library production protocol used. C.A. and Y.H. assisted with
298 sample preparation, processing and tracking. P.J.C., I.M. and M.R.S. oversaw the
299 analysis of mutational signatures and selection analyses. P.J.C., M.H. and S.F.B.
300 wrote the manuscript, with contributions from all authors.

301

302 **COMPETING INTERESTS**

303 The authors declare no competing interests.

304

305 **Correspondence and requests for materials** should be addressed to P.J.C. and
306 M.H.

307

308 **REFERENCES**

- 309 1. The Cancer Genome Atlas Research Network. Comprehensive and
310 Integrative Genomic Characterization of Hepatocellular Carcinoma. *Cell* 169,
311 1327–1341 (2017).
- 312 2. Schulze, K. et al. Exome sequencing of hepatocellular carcinomas
313 identifies new mutational signatures and potential therapeutic targets. *Nat.*
314 *Genet.* 47, 505–511 (2015).

- 315 3. Totoki, Y. et al. Trans-ancestry mutational landscape of hepatocellular
316 carcinoma genomes. *Nat. Genet.* 46, 1267–73 (2014).
- 317 4. Fujimoto, A. et al. Whole-genome sequencing of liver cancers identifies
318 etiological influences on mutation patterns and recurrent mutations in
319 chromatin regulators. *Nat. Genet.* 44, 760–4 (2012).
- 320 5. Letouzé, E. et al. Mutational signatures reveal the dynamic interplay of
321 risk factors and cellular processes during liver tumorigenesis. *Nat. Commun.* 8,
322 (2017).
- 323 6. Kan, Z. et al. Whole-genome sequencing identifies recurrent mutations in
324 hepatocellular carcinoma. *Genome Res.* 23, 1422–1433 (2013).
- 325 7. Guichard, C. et al. Integrated analysis of somatic mutations and focal copy-
326 number changes identifies key genes and pathways in hepatocellular carcinoma.
327 *Nat. Genet.* 44, 694–8 (2012).
- 328 8. Blokzijl, F. et al. Tissue-specific mutation accumulation in human adult
329 stem cells during life. *Nature* 538, 260–264 (2016).
- 330 9. Lodato, M. A. et al. Aging and neurodegeneration are associated with
331 increased mutations in single human neurons. *Science* (80-.). 559, 1–8 (2017).
- 332 10. Lee-Six, H. et al. Population dynamics of normal human blood inferred
333 from somatic mutations. *Nature* 561, 473–478 (2018).
- 334 11. Martincorena, I. et al. High burden and pervasive positive selection of
335 somatic mutations in normal human skin. *Science* (80-.). 348, 880–886 (2015).
- 336 12. Fellous, T. G. et al. Locating the stem cell niche and tracing hepatocyte
337 lineages in human liver. *Hepatology* 49, 1655–63 (2009).
- 338 13. Sigal, S. H. et al. Partial hepatectomy-induced polyploidy attenuates
339 hepatocyte replication and activates cell aging events. *Am. J. Physiol.* 276, G1260-
340 72 (1999).
- 341 14. Stephens, P. J. et al. Massive genomic rearrangement acquired in a single
342 catastrophic event during cancer development. *Cell* 144, 27–40 (2011).
- 343 15. Fernandez-Banet, J. et al. Decoding complex patterns of genomic
344 rearrangement in hepatocellular carcinoma. *Genomics* 103, 189–203 (2014).
- 345 16. Imielinski, M., Guo, G. & Meyerson, M. Insertions and Deletions Target
346 Lineage-Defining Genes in Human Cancers. *Cell* 168, 460-472.e14 (2017).

- 347 17. Martincorena, I. et al. Universal Patterns of Selection in Cancer and
348 Somatic Tissues. *Cell* 171, 1029–1041 (2017).
- 349 18. Torrecilla, S. et al. Trunk mutational events present minimal intra- and
350 inter-tumoral heterogeneity in hepatocellular carcinoma. *J. Hepatol.* 67, 1222–
351 1231 (2017).
- 352 19. Nik-Zainal, S. et al. The life history of 21 breast cancers. *Cell* 149, 994–
353 1007 (2012).
- 354 20. Alexandrov, L. B. et al. Signatures of mutational processes in human
355 cancer. *Nature* 500, 415–421 (2013).
- 356 21. Osorio, F. G. et al. Somatic Mutations Reveal Lineage Relationships and
357 Age-Related Mutagenesis in Human Hematopoiesis. *Cell Rep.* 25, 2308-2316.e4
358 (2018).
- 359 22. Haradhvala, N. J. et al. Mutational Strand Asymmetries in Cancer Genomes
360 Reveal Mechanisms of DNA Damage and Repair. *Cell* 164, 538–549 (2016).
- 361 23. Poon, S. L. et al. Genome-wide mutational signatures of aristolochic acid
362 and its application as a screening tool. *Sci. Transl. Med.* 5, 197ra101 (2013).
- 363 24. Scelo, G. et al. Variation in genomic landscape of clear cell renal cell
364 carcinoma across Europe. *Nat. Commun.* 5, 5135 (2014).
- 365 25. Rushing, B. R. & Selim, M. I. Aflatoxin B1: A review on metabolism,
366 toxicity, occurrence in food, occupational exposure, and detoxification methods.
367 *Food Chem. Toxicol.* 124, 81–100 (2018).
- 368 26. Martincorena, I. et al. Somatic mutant clones colonize the human
369 esophagus with age. *Science* (80-.). 917, 911–917 (2018).
- 370 27. Yokoyama, A. et al. Age-related remodelling of oesophageal epithelia by
371 mutated cancer drivers. *Nature* 1 (2019). doi:10.1038/s41586-018-0811-x
- 372 28. Nault, J. C. et al. Telomerase reverse transcriptase promoter mutation is
373 an early somatic genetic alteration in the transformation of premalignant
374 nodules in hepatocellular carcinoma on cirrhosis. *Hepatology* 60, 1983–92
375 (2014).
- 376 29. Kim, S. K. et al. Comprehensive analysis of genetic aberrations linked to
377 tumorigenesis in regenerative nodules of liver cirrhosis. *J. Gastroenterol.* (2019).
378 doi:10.1007/s00535-019-01555-z

379 30. Zhu, M. et al. Somatic Mutations Increase Hepatic Clonal Fitness and
380 Regeneration in Chronic Liver Disease. *Cell* 1–14 (2019).
381 doi:10.1016/j.cell.2019.03.026
382

383 **FIGURE LEGENDS**

384

385 **Figure 1: Mutational burden observed in non-cancerous hepatocytes.**

386 (A) Burden of SNVs corrected by sensitivity of mutation detection. Each boxplot
387 represents a patient (n=14 patients; 482 microdissections), each dot represents
388 one laser-capture microdissected sample. The grey-to-black intensity of the
389 points reflects the median variant allele fraction (vaf) of mutations in each
390 microdissection. Boxes in the box-and-whisker plots indicate median and
391 interquartile range; whiskers denote range.

392 (B) Burden of insertion-deletion (INDEL) variants (n=14 patients; 482
393 microdissections).

394 (C) Burden of copy number variants (CNVs) and structural variants (SVs),
395 represented as number of unique events per patient.

396 (D) Chromothripsis involving chromosomes 16 and 21 observed in patient
397 PD37111. Black points represent corrected read-depth along the chromosome.
398 Lines and arcs represent structural variants, coloured by orientation of joined
399 ends (purple, tail-to-tail inverted; orange, head-to-head inverted; pale blue,
400 tandem duplication-type orientation; pale green, deletion-type orientation).

401 (E) Chromothripsis involving chromosomes 1 and 3 observed in patient
402 PD37105.

403 (F) Chromothripsis involving chromosomes 2, 5 and 6 observed in patient
404 PD37105 (in a separate clone to panel E).

405

406 **Figure 2: Phylogenetic reconstruction of hepatocyte clones.**

407 (A) Phylogenetic tree constructed from clustering of mutations across
408 microdissected samples in a normal patient (PD36715). Lengths of branches (x
409 axis) indicate numbers of mutations assigned to that branch. Solid lines: nesting
410 is in accordance with the pigeon-hole principle. Dashed lines: nesting is in
411 accordance with the pigeon-hole principle assuming hepatocytes represent 70%
412 of cells. Dotted lines: nesting is only based on clustering, assigning a clone as
413 nested if variant allele fractions of constituent microdissections are lower than
414 those in the parental clone.

415 (B) Representation of branches from the phylogenetic tree in panel A according
416 to their physical coordinates, overlaid onto an H+E stained section. Black points
417 represent branches of the tree sharing no mutations with any other samples;
418 coloured points represent branches with shared clonal relationships (n=26
419 microdissections).

420 (C, D) A second normal liver sample (PD36713; n=30 microdissections).

421 (E, F) Patient with ARLD (PD37105; n=31 microdissections)

422 (G, H) Patient with ARLD (PD37110; n=22 microdissections)

423 (I, J) Patient with NAFLD (PD37114; n=41 microdissections)

424 (K, L) Patient with NAFLD (PD37115; n=34 microdissections)

425 (M, N) Patient with NAFLD (PD37116; 43 microdissections)

426 (O, P) Patient with NAFLD (PD37118; 26 microdissections)

427

428 **Figure 3: Mutational signatures in normal liver, cirrhotic liver and HCC.**

429 (A) Number of somatic substitutions (SNVs; sensitivity-corrected for non-
430 cancerous samples) and insertion-deletion events (INDELs) in each non-cancer
431 microdissection sample (blue points) and associated synchronous HCC (red
432 diamonds).

433 (B) Stacked bar blot showing estimated proportional contributions of each
434 mutational signature to each phylogenetically defined cluster of somatic
435 substitutions. Data generated using a Bayesian hierarchical Dirichlet process.

436 (C) Stacked bar blot showing proportional contributions of signatures in patients
437 with ARLD.

438 (D) Stacked bar blot showing estimated proportional contributions of signatures
439 in patients with NAFLD.

440 (E) Stacked bar blot showing estimated proportional contributions of signatures
441 to 54 cases of HCC from TCGA¹.

442 (F) Number of SNVs attributed to prevalent mutation signatures in each non-
443 cancer microdissection sample (blue circles) and synchronous HCCs (red
444 diamonds). Contributions for the TCGA samples are shown on the right. The y-
445 axis is on a logarithmic scale.

446

447 **Figure 4: The liver as a witness for mutagenic insults occurring throughout**
448 **life.**

449 **(A) Left panel:** Phylogenetic tree of clones in patient PD37111, with each branch
450 coloured by the proportion of mutations in that branch assigned to the different
451 mutational signatures.

452 **Middle panel:** Overlay of the clones represented in (A) onto an H+E stained liver
453 section of patient PD37111 (n=39 microdissections). Colouring of clones is
454 according to the proportion of mutations attributed to Sig. 4, linked to tobacco
455 exposure (blue: low activity of Sig. 4, red: high activity of Sig. 4).

456 **Right panel:** Representative mutation spectrum for samples with low (top) or
457 high (bottom) burden of Sig. 4. The six substitution types are labelled across the
458 top. Within each substitution type, the contribution from the trinucleotide
459 context are shown as 16 bars. The 16 bars are divided into four sets of four bars,
460 grouped by whether an A, C, G or T respectively is 5' to the mutated base, and
461 within each group of four by whether A, C, G or T is 3' to the mutated base.

462 **(B)** Overlay of mutational signatures onto phylogenetic tree of clones in patient
463 PD37107 (n=41 microdissections). Colouring of clones in the middle panel is
464 according to Sig. 22, linked to the aristolochic acid carcinogen.

465 **(C)** Overlay of mutational signatures onto phylogenetic tree of clones in patient
466 PD36714 (n=35 microdissections). Colouring of clones in middle panel is
467 according to Sig. 24, linked to the carcinogen aflatoxin-B₁.

468 **(D)** Overlay of mutational signatures onto phylogenetic tree of clones in patient
469 PD37113 (n=37 microdissections). Cluster 10 has many mutations attributed to
470 Sig. 9, linked to the somatic hypermutation process in B lymphocytes.

471

472 **EXTENDED FIGURE LEGENDS**

473

474 **Extended Figure 1: Sensitivity analysis of SNV calls.**

475 (A) Overview schematic of the experimental and analytical approach.

476 (B) Examples of the variant allele fractions (VAFs) of variants from unrelated
477 (top) and related (bottom) microdissection sample pairs from four donors (left
478 to right). X-axis represents the VAF of sample 1 from each pair; Y-axis represents
479 the VAF of sample 2. Each dot represents one variant. Red: variants called in both
480 samples, yellow: variants called in sample 1, blue: variants called in sample 2.

481 (C) Histogram of sensitivities calculated for each sample pair.

482 (D) Heatmap of modelled sensitivity at different values of VAF and coverage.
483 Overlaid dots represent sample pairs used to fit model.

484 (E) Relationship of VAF, sensitivity and coverage according to fitted model of
485 sensitivity. Overlaid dots represent sample pairs used to fit model.

486 (F) Comparison of calculated (x-axis) and fitted (y-axis) sensitivity for each
487 sample pair (n=34 pairs of samples). The R² value quoted is a Pearson's
488 correlation coefficient.

489 (G) Proportion of hepatocytes that are multinucleated in samples analysed here,
490 estimated by counting 500 cells in each H&E section (n=14 patients). Each point
491 represents the proportion of a patient in the study. The horizontal bars
492 represent the mean for that aetiological group.

493

494 **Extended Figure 2: Copy number and structural variants in chronic liver**

495 **disease. (A, B)** Genome-wide copy number profiles for two samples. Black
496 points represent read-depth of discrete windows along the chromosome,
497 corrected to show overall copy number. Arm-level and whole chromosome gains
498 and losses are evident.

499 (C-H) Focal copy number changes and structural variants. Black points represent
500 read-depth of discrete windows along the chromosome, corrected to show
501 overall copy number. Lines and arcs represent individual structural variants,
502 coloured by the orientation of the joined ends (purple, tail-to-tail inverted;
503 orange, head-to-head inverted; pale blue, tandem duplication-type orientation;

504 pale green, deletion-type orientation). Events affecting known HCC genes are
505 marked with labelled arrows (panels C, E, F).

506

507 **Extended Figure 3: Events affecting known HCC genes in cohort.**

508 (A) Distribution of somatic point mutations in individual microdissections (x
509 axis) affecting known HCC genes (y axis). The inset to the left shows the
510 frequency of events in individual genes. The inset to the bottom shows the
511 aetiology attributed to the sample, and whether the sample was drawn from
512 non-cancerous hepatocytes (left) or HCC (right).

513 (B) Genomic position of single nucleotide substitutions (SNVs; light blue strip,
514 top) and insertion-deletions (INDELs; dark blue strip, bottom) detected in *ALB*,
515 the gene encoding albumin.

516 (C) Relationship of gene expression in liver tissue (x axis) and proportion of
517 indels as a fraction of all point mutations (y axis). The grey line represents a
518 Poisson regression model with a significant (two-sided likelihood ratio test; $p <$
519 10^{-16}) coefficient for gene expression as a predictor for the ratio of indels
520 ($n=5458$ genes included in model). The grey ribbon represents the 99%
521 confidence interval of the parameter estimates.

522

523 **Extended Figure 4: Phylogenetic reconstruction of hepatocyte clones in**
524 **non-cirrhotic liver samples.**

525 Left column: Heatmap representing the clustering of the variants observed in
526 each microdissection sample (x-axis) of the non-cirrhotic livers. Each cluster (y-
527 axis) contains mutations for which variant allele fractions across samples are
528 very similar. The colour scale of the boxes represents the estimated mean variant
529 allele fraction for that cluster in that sample.

530 Middle column: Phylogenetic trees constructed from the clustering information.
531 Solid lines: nesting is in accordance with the pigeon-hole principle. Dashed lines:
532 nesting is in accordance with the pigeon-hole principle assuming the pool of
533 hepatocytes to be 70% of cells. Dotted lines: nesting is only based on clustering,
534 assigning a clone as nested if its constituent LCMs are a subset of LCMs in the
535 parental clone. Details given in Supplementary Methods.

536 Right column: Representation of clones according to the physical coordinates of
537 the LCM samples, overlaid onto H&E stained sections (top), with Masson's
538 trichrome and Oil Red-O sections also shown (bottom). Locations of
539 immune/inflammatory cell infiltrates are marked with yellow rings. Sample sizes
540 were for PD36713, n=30 microdissections; PD36714, n=35 microdissections;
541 PD36715, n=26 microdissections; PD36717, n=42 microdissections; PD36718,
542 n=32 microdissections.

543

544 **Extended Figure 5: Phylogenetic reconstruction of hepatocyte clones in**
545 **alcohol-related cirrhosis.**

546 Analogous to Extended Figure 4, representing the cirrhotic livers of donors
547 PD37105, PD37107, PD37110 and PD37111. The pictures in the right column
548 are of H&E stains on the top, with Masson's trichrome and a macroscopic
549 photograph of the liver on the bottom, with HCCs indicated by arrows. Locations
550 of immune/inflammatory cell infiltrates are marked with yellow rings. Sample
551 sizes were for PD37105, n=31 microdissections; PD37107, n=41
552 microdissections; PD37110, n=22 microdissections; PD37111, n=39
553 microdissections.

554

555 **Extended Figure 6: Phylogenetic reconstruction of hepatocyte clones in**
556 **non-alcoholic fatty liver disease with cirrhosis.**

557 Analogous to Extended Figure 4, representing the cirrhotic livers of donors
558 PD37113, PD37114, PD37115, PD37116 and PD37118. The pictures in the right
559 column are of H&E stains on the top, with Masson's trichrome and a macroscopic
560 photograph of the liver on the bottom, with HCCs indicated by arrows. Locations
561 of immune/inflammatory cell infiltrates are marked with yellow rings. Sample
562 sizes were for PD37113, n=37 microdissections; PD37114, n=41
563 microdissections; PD37115, n=34 microdissections; PD37116, n=43
564 microdissections; PD37118, n=26 microdissections.

565

566 **Extended Figure 7: Mutation spectrum of individual microdissections**

567 From each donor, we chose 5 clones to represent the heterogeneity in
568 trinucleotide context mutation spectra. The six substitution types are shown in

569 the panel across the top of each clone's data. Within each panel, the contribution
570 from the trinucleotide context (bases immediately 5' and 3' of the mutated base)
571 are shown.

572

573 **Extended Figure 8: Details of mutational signature extractions**

574 (A) Dot plots showing the concordance for signature attributions between the
575 two signature algorithms (n=479 microdissections). Mutational signatures on
576 the y axis were extracted using non-negative matrix factorisation and on the x
577 axis using a Bayesian hierarchical Dirichlet process. Quoted R values are
578 Pearson's correlation coefficients.

579 (B) Signatures extracted by non-negative matrix factorisation. The six
580 substitution types are shown in the panel across the top of each clone's data.
581 Within each panel, the contribution from the trinucleotide context (bases
582 immediately 5' and 3' of the mutated base) are shown.

583 (C) Signatures extracted by the Bayesian hierarchical Dirichlet process, as for
584 panel B. Where a signature matches one from panel B, it is shown on the same
585 row.

586

587 **Extended Figure 9: Transcription strand bias in mutational patterns**

588 (A) Transcription strand bias of T>C mutations at A[T]D context before and after
589 transcription start sites of highly expressed liver genes.

590 (B) Bar plots representing the numbers of C>A variants on the transcribed and
591 non-transcribed strand. Each hepatocyte clone is represented individually (x-
592 axis). Note the strand bias in the highly mutated clones of PD37111, where the
593 tobacco signature is most active – the strand bias indicates the damaged base is
594 the guanine, as expected for polycyclic aromatic hydrocarbons.

595 (C) Bar plots representing the numbers of T>A variants on the transcribed and
596 non-transcribed strand. Each hepatocyte clone is represented individually (x-
597 axis). Note the strand bias in the highly mutated clones of PD37107, where the
598 aristolochic acid signature is most active – the strand bias indicates the damaged
599 base is the adenine, as expected for polycyclic aromatic hydrocarbons.

600

601 **Extended Figure 10: Mutations in a B lymphocyte clone in a cirrhotic liver**

602 (A) Illustration of a portion of the B-cell receptor (*IGH*) region on chromosome
603 14. Shown are the coverage tracks of an LCM sample that does not belong to the
604 lymphocyte lineage (top) and a sample that belongs to the lymphocyte lineage
605 (middle). In the center of the displayed region there is a drop of copy number in
606 the lymphocyte track, indicating a structural rearrangement. The bottom track
607 shows the paired-end reads that contribute to a rearrangement event in the
608 lymphocyte sample, co-localised with the drop in copy number.

609 (B) Application of the pigeonhole principle – if two clusters of heterozygous
610 mutations in regions of diploid copy number are in different cells, then their
611 median variant allele fractions must sum to ≤ 0.5 (if they sum to >0.5 , equivalent
612 to a combined cellular fraction of >1 , there must be some cells that carry both
613 sets of mutations – hence one cluster would have a subclonal relationship with
614 the other). Cluster 10 is the cluster with the unique VDJ rearrangement of *IGH*
615 shown in panel A and the large number of mutations attributed to signature 9.
616 Clearly, samples from clusters 2, 11 and 55 etc have VAFs which, when combined
617 with cluster 10, sum to >0.5 . Therefore, they must be subclonal to cluster 10,
618 even though they do show signature 9.

619 (C-H) Representative pairwise decision graphs for clusters of mutations. Median
620 cellular fraction is shown for pairs of clusters across every sample from the
621 patient. Where at least one sample falls above / to the right of the $x+y=1$ diagonal
622 line, those two clusters must share a nested clonal-subclonal relationship.

623

624 **Methods**

625

626 **SAMPLES AND SEQUENCING**

627 **Samples**

628 Patients recruited at Addenbrooke's Hospital, Cambridge gave written informed
629 consent with approval of the Local Research Ethics Committee (16/NI/0196).

630

631 Normal liver samples were obtained from patients with liver metastases from
632 colorectal carcinoma (CRC). The liver specimens were obtained from resected
633 liver distal to the metastases, that were confirmed on histology. None of the
634 patients had undergone neo-adjuvant systemic therapy; one patient had
635 undergone pre-operative portal vein embolisation (PD36718) to the ipsilateral
636 liver lobe. Liver tissue from patients with chronic liver disease (CLD) was
637 derived from explanted diseased livers at the time of transplantation. All of the
638 patients were identified as having ARLD or NAFLD by clinical history to the
639 transplant hepatology and addiction psychiatry teams, as well as explanted liver
640 histology. None of the patients had undergone trans-arterial chemo-embolisation
641 (TACE) or other locoregional therapy on the transplant waiting list, except
642 PD37118 who underwent a single treatment to their HCC with TACE. All of the
643 CLD patients, except one (PD37105), demonstrated significant pre-operative
644 impairment of liver function as evidenced by a UKELD of >50.

645

646 The explant liver histology was reviewed by a specialist liver histopathologist
647 (SED), blinded to the sequencing results. The normal liver specimens had no
648 fibrosis and no evidence of chronic liver disease; the explanted diseased livers
649 uniformly demonstrated cirrhosis and HCC. The background liver histology was
650 scored according to the Kleiner system³¹ on FFPE samples away from the HCC
651 and the fresh frozen block used for the sequencing analysis. The Kleiner score
652 assesses the presence of steatosis, lobular inflammation and hepatocyte
653 ballooning to generate a cumulative NAS score. The presence or absence of
654 cellular or nodular dysplasia was globally assessed in clinical FFPE samples
655 (Supplementary table 1), as well as specifically assessed in the fresh-frozen block
656 used for the laser capture microdissection and sequencing (Supplementary table

657 1). Serial H&E-stained sections from the frozen block did not demonstrate
658 dysplasia in any of the cases (Supplementary table 1). Further, there was no
659 evidence of CRC or HCC on histological review of the fresh-frozen block used for
660 sequencing.

661

662 All tissue samples were snap-frozen in liquid nitrogen and stored at -80°C in the
663 Human Research Tissue Bank of the Cambridge University Hospitals NHS
664 Foundation Trust.

665

666 **Preparation of tissue sections**

667 Tissue biopsies were embedded in Optimal Cooling Temperature (OCT,
668 ThermoFisher) medium at -25°C. Sections were cut at a thickness of 20µm using
669 a Leica Cryotome and transferred onto PEN membrane slides (ThermoFisher).
670 For fixation, slides were treated with 70% ethanol at room-temperature for
671 2min. Slides were washed twice in 10% phosphate buffered saline (PBS) at
672 room-temperature for 10s. For staining, slides were incubated in haematoxylin
673 for 10s and rinsed twice in water. Slides were then incubated in eosin for 5s and
674 rinsed once in water. Slides were washed twice with 70% ethanol for 5s, twice
675 with 100% ethanol for 5s, and in xylene for 5s. Storage was at -20°C. Additional
676 sections were stained for H&E, Masson's Trichrome and Oil Red O by standard
677 laboratory techniques. All slides were scanned on a Leica AT2 at ×20
678 magnification and a resolution of 0.5µm per pixel.

679

680 **Laser Capture Microdissection (LCM)**

681 Microdissection was performed using a LCM (Leica Microsystems LMD 7000).
682 For each biopsy, 48 microdissections were cut with a target size of 20,000µm²,
683 corresponding to about 400 hepatocyte cells. Images were taken before and after
684 LCM.

685

686 **Sample lysis and DNA preparation**

687 LCM biopsies were lysed using the Arcturus PicoPure DNA Extraction Kit
688 (ThermoFisher) following the manufacturer's instructions. DNA libraries for

689 Illumina sequencing were prepared using a protocol optimized for low input
690 amounts of DNA, as described³².

691

692 **Whole-genome sequencing**

693 Paired-end sequencing reads (150bp) were generated using the Illumina X10
694 platform for 400 samples, resulting in a target coverage of 30x-70x per sample.
695 To avoid the known index-hopping artefact, we chose to avoid multiplexing
696 samples and instead sequenced one sample per flow cell lane. To increase
697 coverage for a subset of 96 samples, we used multiplexing and achieved 70x
698 coverage. In addition to the LCM samples we also sequenced a bulk sample for
699 each biopsy and (where available) associated hepatocellular carcinoma (HCC).

700

701 The healthy liver samples came from wide resections of hepatic metastases of
702 colorectal cancer. In each case, we sequenced the metastasis – this did not reveal
703 any mutations shared between the colorectal cancer and liver, nor any variants
704 shared by all liver samples absent from the colorectal cancer (beyond regions of
705 loss-of-heterozygosity in the cancer). Likewise, for the cirrhotic liver samples, we
706 sequenced the matched HCC, not revealing sharing of mutations. In one case, we
707 sequenced microdissections of the fibrotic tissue, and here also did not find
708 mutations restricted to all liver cells.

709

710 Sequencing data were mapped to the human genome, GRCh37d5, using the
711 BWA-Mem algorithm.

712

713 **VARIANT CALLING**

714 **SNV calling**

715 Substitution variants were called using the Cancer Variants through Expectation
716 Maximisation (CaVEMan) algorithm³³, using the bulk sample of the liver biopsy
717 as the matched normal. As part of the algorithm, the variants were annotated
718 using VAGrENT³⁴. Variant calls for bulk sequencing data of the cancer samples
719 were not further filtered. For sequencing of LCMs, post-filtering was performed
720 in three steps:

721

722 1. *Removal of duplicate counts:* we noticed instances where variant bases were
723 counted twice due to the overlap of paired-end sequencing reads. We removed
724 such double counting and re-evaluated variant calls after taking double counts
725 into account.

726

727 2. *Removal of variants introduced during library preparation:* we noticed the
728 presence of variants introduced due to incorrect processing of cruciform DNA.
729 Erroneous variants were often present in inverted repeats and frequently
730 accompanied by another proximal (~ 1-30bp distance). These inverted repeats
731 can form cruciform DNA prior to DNA isolation or during library preparation.
732 The library preparation protocol employed can incorrectly process these
733 secondary DNA structures and inadvertently introduce one or more erroneous
734 variants. For every variant the standard deviation (SD) and median absolute
735 deviation (MAD) of the variant position within the read was separately
736 calculated for positive and negative strand reads.

737 In the case that the variant was supported by a low number of reads for a
738 particular strand, the filtering was based on the statistics determined from the
739 reads derived from the other strand. It was required that either:

740 1. $\leq 90\%$ of supporting reads report the variant within the first 15% of the
741 read as calculated from the alignment start.

742 2. Or, that the $MAD > 0$ and $SD > 4$.

743

744 In the case that sufficient reads supporting the variant were available for both
745 strands it was required for both strands separately that either:

746 1. $\leq 90\%$ of supporting reads report the variant within the first 15% of the
747 read as calculated from the alignment start.

748 2. Or, that the $MAD > 2$ and $SD > 2$.

749 3. Or, that at least one strand has fulfills the criteria $MAD > 1$ and $SD > 10$.

750

751 3. *Comparison with an independent panel:* to remove variant calls at badly-
752 mapping sites, we compared variant calls in the sequenced samples of each
753 donor biopsy with samples from all unrelated donors in our cohort. For each
754 variant site we expected the reference base to be dominant and conversely

755 expected badly-mapping sites to contain frequent non-reference base counts.
756 Thus, we counted the numbers of A, C, G, T, insertion and deletion calls at each
757 variant site across all unrelated samples, resulting in a large “pileup” table. The
758 dominance of the reference base was evaluated at each variant site using the
759 entropy purity metric E :

$$E = - \sum_{i \in \{A, C, G, T, Ins, Del\}} P(x_i) \ln P(x_i)$$

760 where x is the count of base i and the $P(x_i)$ are the fractions of base calls. Values
761 of E close to 0 indicate that almost all reads in the independent panel contain a
762 single base. Higher values of E indicate a mix of base calls at the site. To identify
763 an optimal threshold of E for the filtering of variant sites, we evaluated the
764 entropy metric against a labelled dataset of variant calls. Specifically, during the
765 clustering of variants using the Bayesian Dirichlet process (described below), we
766 identified clusters that had variants with low allele frequency present in all
767 dissections from the same donor. Manual inspection showed that such variants
768 occurred at badly-mapping sites. Thus, we labelled variant sites in those clusters
769 as “badly-mapping” and were able to use the Area-Under-the-Receiver-Operator-
770 Curve to identify a threshold value E_{Thr} of 0.16 that allowed to separate the two
771 labelled variant groups with an AUC of 0.99.

772

773 **Bayesian Dirichlet process for clustering VAFs across multiple samples**

774 We extend the model previously developed for clustering variant allele fractions
775 (VAFs) of mutations called in a single sample¹⁹ to mutation data across multiple
776 samples from the same individual. In normal somatic cells, the vast majority of
777 the genome retains its normal, diploid copy number, which means that we can
778 cluster the VAFs directly (excluding mutations on the X and Y chromosomes in
779 males) – this has the considerable advantage that the Dirichlet Process model we
780 build can rely directly on conjugate prior distributions. The model includes a
781 potential split-merge step at each cycle of the Gibbs sampler, following a
782 previously described Metropolis-Hastings proposal for conjugate distributions³⁵.
783 The algorithm could be extended to include a correction for different copy
784 number states in given samples for a particular mutation through, for example, a
785 Metropolis-Hastings update, but at considerable computational cost. The full

786 mathematical development of the model is detailed in the **Supplementary**
787 **Methods**.

788

789 We ran the Gibbs sampler for 15,000 iterations, dropping the first 10,000 as a
790 burn-in. We used the ECR algorithm³⁶, implemented in the R package
791 label.switching, to resolve the label switching problem associated with mixture
792 models. We dropped clusters containing <100 variant sites.

793

794 **Phylogenetic tree construction**

795 Phylogenetic trees were constructed manually using the pigeonhole principle as
796 described previously¹⁹. In short, each cluster identified using the Bayesian
797 Dirichlet process represented a branch of the phylogenetic tree. Nesting of trees
798 was identified with three different levels of certainty, illustrated on a pair of
799 branches A and B:

- 800 1. In case the median VAFs of A and B exceeded 100%, the pigeonhole
801 principle defines that A and B are nested.
- 802 2. We can assume that non-hepatocyte cells constitute a sizeable fraction of
803 each LCM sample. Assuming a non-hepatocyte fraction of 30% we nested
804 branches when VAFs of A and B exceeded 70%. This non-hepatocyte
805 fraction was chosen as a conservative estimate of the fraction of cells
806 intermixed in our microdissections that are not derived from the
807 hepatocyte clone, based on observed VAF peaks in our data together with
808 single-cell RNA sequencing data from liver tissue.
- 809 3. If identical LCMs are members of both A and B, it is highly likely that A
810 and B are nested, rather than independent branches. Thus, we also nested
811 branches where the LCMs in one branch were a subset of the LCMs in the
812 other (parental) branch.

813

814 In each nesting scenario, we defined the parental branch to be the one with the
815 higher median VAF in the contained LCMs. We highlighted the evidence level for
816 nesting in each representation of phylogenetic trees, marking branches with
817 evidence level 1 with a solid line, level 2 with a dashed line and level 3 with a
818 dotted line.

819

820 **Analysis of driver variants**

821 We curated a list of genes that have been found to be significantly mutated in
822 liver cancers in a selection of published studies^{1-4,6,7,37-39}, as represented in
823 Supplementary Table 5. Using the VAGrENT annotations³⁴, we counted any
824 regulatory, missense, nonsense, frameshift or essential splice variant as a
825 potential driver variant. To systematically identify genes under mutagenic
826 selection, we used the dN/dS method¹⁷ that screens for genes with an excess of
827 non-synonymous mutations compared to that expected from the synonymous
828 mutation rate.

829

830 **Sensitivity correction**

831 We identified 138 pairs of LCMs with a midpoint-to-midpoint distance of <
832 500µm and at least one shared cluster according to the Bayesian Dirichlet
833 process. These LCMs we assumed to represent the same clone, thus providing an
834 opportunity to calculate the sensitivity of calling a variant present in one LCM in
835 the other. If we assume the sensitivity is the same in both samples, then the
836 maximum likelihood estimate for the sensitivity, when mutations not called in
837 either sample are unobserved, is given by:

$$s = \frac{2n_2}{n_1 + 2n_2}$$

838 where n_2 is the number of variants called in both LCMs in each pair and n_1 is the
839 number of variants called only in one of the two LCMs. To evaluate the
840 relationship of sensitivity with depth-of-coverage and VAF, we performed a
841 logistic regression of sensitivity against these two predictors using the `lm()`
842 function of the R programming language. The model fit was then used to
843 calculate sensitivity for any LCM sample, given the coverage and VAF of the
844 sample.

845

846 **Mutation burden analysis**

847 We used a linear mixed effects model to fit the number of variants per LCM
848 sample against each individual's disease aetiology (normal or cirrhotic) and age.
849 We defined the individual's ID as a random effect. The slope of the age coefficient

850 was allowed to vary with the random effect. To facilitate the analysis, we used
851 the lmer() function available from the lme4 package of the R programming
852 language. To determine the significance of the aetiology and age coefficients, we
853 used ANOVA analysis to perform a X^2 test comparing our model with models
854 omitting the aetiology and age coefficients, respectively.

855

856 **Deep targeted sequence validation of mutation calls**

857 For 96 of the microdissections sequenced by whole genome sequencing, we
858 performed a deep targeted sequencing validation using an Agilent RNA bait-set
859 covering 350 recurrently mutated cancer genes. Among these genes, a total of 17
860 mutations were identified in the whole genome sequencing data from the 96
861 samples – of these, 16 (94%) were validated, at comparable variant allele
862 fractions, in the targeted deep sequencing data.

863

864 **INDEL calling**

865 INDELS were called using cgpPindel⁴⁰. Variant calls for bulk sequencing data of
866 the cancer samples were not further filtered. To remove artefactual calls from
867 the LCM-derived data, we performed two post-filtering steps:

868

869 *1) Assignment to SNV-based clusters:* we evaluated how well the VAF distribution
870 of each INDEL across the LCMs from the same donor compared with the VAF
871 distribution of each SNV-based cluster as identified by the Bayesian Dirichlet
872 process. Given an INDEL in one LCM sample, we thus counted its occurrence in
873 all related LCMs and assigned the resulting VAF profile to the SNV clusters' VAF
874 profiles using a Bayes' classifier. We noticed that many INDELS were assigned to
875 SNV clusters with <100 variants, which we had previously removed from the
876 SNV analysis. On closer inspection we noticed that those INDELS had low VAF
877 and occurred frequently in badly-mapping regions. We thus discarded INDELS
878 assigned to those clusters.

879

880 *2) Filtering based on beta-binomial overdispersion parameter:* we noticed that
881 many INDELS occurred with low VAF in a large number of LCMs from the same
882 donor and were, thus, likely to be artefactual. To systematically identify such

883 INDELS, we fitted the beta-binomial distribution to the variant counts of each
884 INDEL across the LCMs from the same donor. Fitted parameter ρ , the
885 overdispersion parameter, was used to filter INDEL calls. A high value for
886 parameter ρ (overdispersion) occurs when some LCMs have many variant read
887 counts and others few or none. Conversely, a low value occurs when all LCMs
888 have a similar number of variant counts (no overdispersion). Based on manual
889 inspection, we removed variant calls with $\rho < 0.02$.

890

891 **Copy number calling**

892 CNs were called using the ASCAT algorithm⁴¹, assuming an expected ploidy of 4
893 (to allow for physiologically polyploid hepatocytes) and 60% non-hepatocyte cell
894 contamination for all samples. Robustness testing around these starting points
895 (different expected ploidy or purity values) found that the specific values used
896 did not materially affect the output. Variant calls for bulk sequencing data of the
897 cancer samples were not further filtered. To remove artefactual variants from
898 the LCM-derived data, we employed the SNV-based phylogenetic information.
899 The genome was segmented into 500bp bins and the ASCAT-based copy number
900 of each bin was calculated. Using the binned CN data we calculated the median
901 CN in each LCM sample and ASCAT event. For each ASCAT event and LCM sample
902 we assigned its absolute deviation from the diploid state. We compared each
903 ASCAT event's CN profile across the LCM samples with the VAF profile of each
904 SNV cluster using cosine similarity (described below) to identify the most similar
905 SNV cluster. Within each SNV cluster we proceeded to merge overlapping ASCAT
906 events. Using manual inspection, we decided to keep ASCAT events if they 1) had
907 a cosine similarity of < 0.1 to an SNV cluster and 2) if their assigned SNV cluster
908 was not removed during SNV analysis due to having < 100 assigned SNVs.

909

910 **Structural variant calling**

911 SVs were called using the BRASS algorithm⁴²
912 (<https://github.com/cancerit/BRASS>). Variant calls for bulk sequencing data of
913 the cancer samples were not further filtered. To remove artefactual variants
914 from the LCM-derived data, we employed post-processing filters. Manual
915 inspection of the sequencing reads identified for each SV showed that many

916 reads were identical except for frame-shifts at repetitive sites. We decided that
917 such reads represented duplicates and designed a filter to systematically remove
918 these. We removed SVs supported by <2 reads after duplicate removal. Each
919 remaining SV call was manually inspected.

920

921 **Clone size calculation**

922 We determined the midpoint coordinates of each LCM manually from the
923 microscopy images collected during dissection. For each LCM belonging to a
924 clone as determined by the Bayesian Dirichlet process, we used the function
925 *chull* of the R programming language to identify the coordinates of the convex
926 hull that included all LCMs. We identified the midpoint of each polygon as the
927 average coordinate of all convex hull vertices. The size of the clone was then
928 assigned to be the Euclidean distance between each convex hull vertex and the
929 polygon's midpoint. For clones that only consisted of a single LCM, we assigned
930 the minimum clone size discovered across all clones.

931

932 **Extraction of mutational signatures from SNV contexts using HDP**

933 Mutational signatures were extracted using the HDP package
934 (<https://github.com/nicolaroberts/hdp>) relying on the hierarchical Bayesian
935 Dirichlet process. The units of signature extraction were mutations assigned to
936 individual branches of the phylogenetic tree, grouped per patient, from the LCM
937 data. In addition, to provide a comparison against signatures extracted in HCCs,
938 we added catalogues of somatic substitutions from 54 whole genomes sequenced
939 by the TGCA, analysed using the same core algorithms as used for the LCM data.
940 The tool was used without defining prior signatures. As hyperparameters we set
941 alpha and beta to 6 for the alpha clustering parameter. Extraction was started
942 with 40 data clusters (parameter 'initcc'). The Gibbs sampler was run with
943 10,000 burn-in iterations (parameter 'burnin'). With a spacing of 50 iterations
944 (parameter 'space'), 50 iterations were collected (parameter 'n'). After each
945 Gibbs sampling iteration, 3 iterations of concentration parameter sampling were
946 performed (parameter 'cpiter'). Resulting signatures were compared to
947 published signatures^{20,43} using the cosine similarity metric described below.
948 Extracted signatures with cosine similarity >0.9 compared to a known signature

949 from either the COSMIC²⁰ or PCAWG⁴³ catalogue of signatures were assigned the
950 name of the known signature with the highest similarity. Extracted signatures
951 with cosine similarity <0.9 to any of the known signatures were assigned new
952 names, indexed with letters A, B, and C.

953

954 **Extraction of mutational signatures from SNV contexts using SigProfiler**

955 We used SigProfiler to extract mutational signatures, relying on the non-negative
956 matrix factorization (NNMF) method⁴⁴. In particular, we report the “Decomposed
957 Solution” output by the package.

958

959 **Cosine similarity calculation**

960 To compare two vectors A and B, cosine similarity was calculated as follows:

$$similarity = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

961

962 **Analysis of INDEL proportion and gene expression**

963 A list of transcribed regions was retrieved from ENSEMBL using the BioMart
964 package⁴⁵. We identified the subset of INDEL and SNV variants that overlapped
965 with the transcribed regions. The proportion of INDELS in comparison to the
966 total number of INDELS and SNVs per gene was calculated. Gene expression was
967 assigned using the “liver” dataset from the Genotype-Tissue Expression project⁴⁶.
968 To test for the relationship of gene expression on INDEL proportion, we fit a
969 Poisson regression using the *glm* function of the R programming language. We
970 modelled the number of INDELS per gene against an offset of the total number of
971 variants per gene and the gene’s expression.

972

973 **Analysis of T>C transcription strand bias at transcription start sites**

974 We performed this analysis analogously to a published approach²². In short, we
975 retrieved the genomic coordinates of transcription start sites of the all
976 overexpressed genes in the liver (GTEx⁴⁶). We tiled the 10 kilobases up- and
977 downstream of the transcription start site into 1,000bp bins. We overlapped all
978 T>C (transcribed) and A>G (untranscribed) variant calls with the tiled regions

979 and summed the number of variants in each tile across all included genes. We
980 also extracted the number of T and A bases in each tile. To test whether strand
981 bias was significant only in transcribed regions, we fit a Poisson regression for
982 the number of variant calls against the following predictors: strand (transcribed
983 / untranscribed), distance from TSS (0 for upstream, 1 for downstream),
984 aetiology (cirrhosis, no cirrhosis) and used the number of T and A bases in each
985 tile as the offset variable.

986

987 **Analysis of C>A and T>A transcription strand bias**

988 We used the MutationalPatterns package⁴⁷ to assign the transcription state for
989 each C>A variant. We retrieved the genomic coordinates of all transcribed
990 regions from ENSEMBL using the BioMaRt package⁴⁵ and extracted the
991 frequencies of C and G nucleotides in these regions. To test for significance of
992 transcription strand bias, we performed a Poisson regression for the number of
993 C>A variants in each sample and transcription strand against factor variables for
994 the transcription strand, the patient ID and an interaction term for the two
995 factors. We used the C, G nucleotide frequency as an offset variable. To test for
996 significance of transcription strand bias for a given donor, we coded the patient
997 ID in a binary fashion: “1” for the target donor, “0” otherwise. We proceeded
998 analogously to test for transcription strand bias of T>A variants, using A and T
999 nucleotide frequencies as the offset.

1000

1001 **REFERENCES (Continued from main text references)**

- 1002 31. Kleiner, D. E. et al. Design and validation of a histological scoring system
1003 for nonalcoholic fatty liver disease. *Hepatology* 41, 1313–1321 (2005).
- 1004 32. Lee-Six, H. et al. The landscape of somatic mutation in normal colorectal
1005 epithelial cells. *bioRxiv* 416800 (2018). doi:10.1101/416800
- 1006 33. Jones, D. et al. cgpCaVEManWrapper: Simple Execution of CaVEMan in
1007 Order to Detect Somatic Single Nucleotide Variants in NGS Data. in *Current*
1008 *Protocols in Bioinformatics* 2016, 15.10.1-15.10.18 (John Wiley & Sons, Inc.,
1009 2016).
- 1010 34. Menzies, A. et al. VAGrENT: Variation Annotation Generator. *Curr. Protoc.*
1011 *Bioinformatics* 52, 15.8.1-15.8.11 (2015).

- 1012 35. Dahl, D. B. An improved merge-split sampler for conjugate Dirichlet
1013 process mixture models. Univ. Wisconsin-Madison Tech. Rep. 1086, 1–32 (2003).
- 1014 36. Papastamoulis, P. label.switching: An R Package for Dealing with the Label
1015 Switching Problem in MCMC Outputs. J. Stat. Softw. 69, Code Snippet 1 (2015).
- 1016 37. Fujimoto, A. et al. Whole-genome mutational landscape of liver cancers
1017 displaying biliary phenotype reveals hepatitis impact and molecular diversity.
1018 Nat. Commun. 6, 1–8 (2015).
- 1019 38. Cleary, S. P. et al. Identification of driver genes in hepatocellular
1020 carcinoma by exome sequencing. Hepatology 58, 1693–702 (2013).
- 1021 39. Ahn, S.-M. et al. Genomic portrait of resectable hepatocellular carcinomas:
1022 implications of RB1 and FGF19 aberrations for patient stratification. Hepatology
1023 60, 1972–82 (2014).
- 1024 40. Raine, K. M. et al. cgppindel: Identifying Somatically Acquired Insertion
1025 and Deletion Events from Paired End Sequencing. Curr. Protoc. Bioinformatics
1026 52, 15.7.1-15.7.12 (2015).
- 1027 41. Raine, K. M. et al. ascatNgs: Identifying Somatically Acquired Copy-
1028 Number Alterations from Whole-Genome Sequencing Data. in Current Protocols
1029 in Bioinformatics 2016, 15.9.1-15.9.17 (John Wiley & Sons, Inc., 2016).
- 1030 42. Campbell, P. J. et al. Identification of somatically acquired rearrangements
1031 in cancer using genome-wide massively parallel paired-end sequencing. Nat.
1032 Genet. 40, 722–9 (2008).
- 1033 43. Alexandrov, L. et al. The Repertoire of Mutational Signatures in Human
1034 Cancer. bioRxiv 322859 (2018). doi:10.1101/322859
- 1035 44. Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M.
1036 R. Deciphering Signatures of Mutational Processes Operative in Human Cancer.
1037 Cell Rep. 3, 246–259 (2013).
- 1038 45. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for
1039 the integration of genomic datasets with the R/Bioconductor package biomaRt.
1040 Nat. Protoc. 4, 1184–91 (2009).
- 1041 46. GTEx Consortium. Genetic effects on gene expression across human
1042 tissues. Nature 550, 204–213 (2017).

1043 47. Blokzijl, F., Janssen, R., van Boxtel, R. & Cuppen, E. Mutational Patterns:
1044 comprehensive genome-wide analysis of mutational processes. *Genome Med.* 10,
1045 33 (2018).
1046

Figure 1

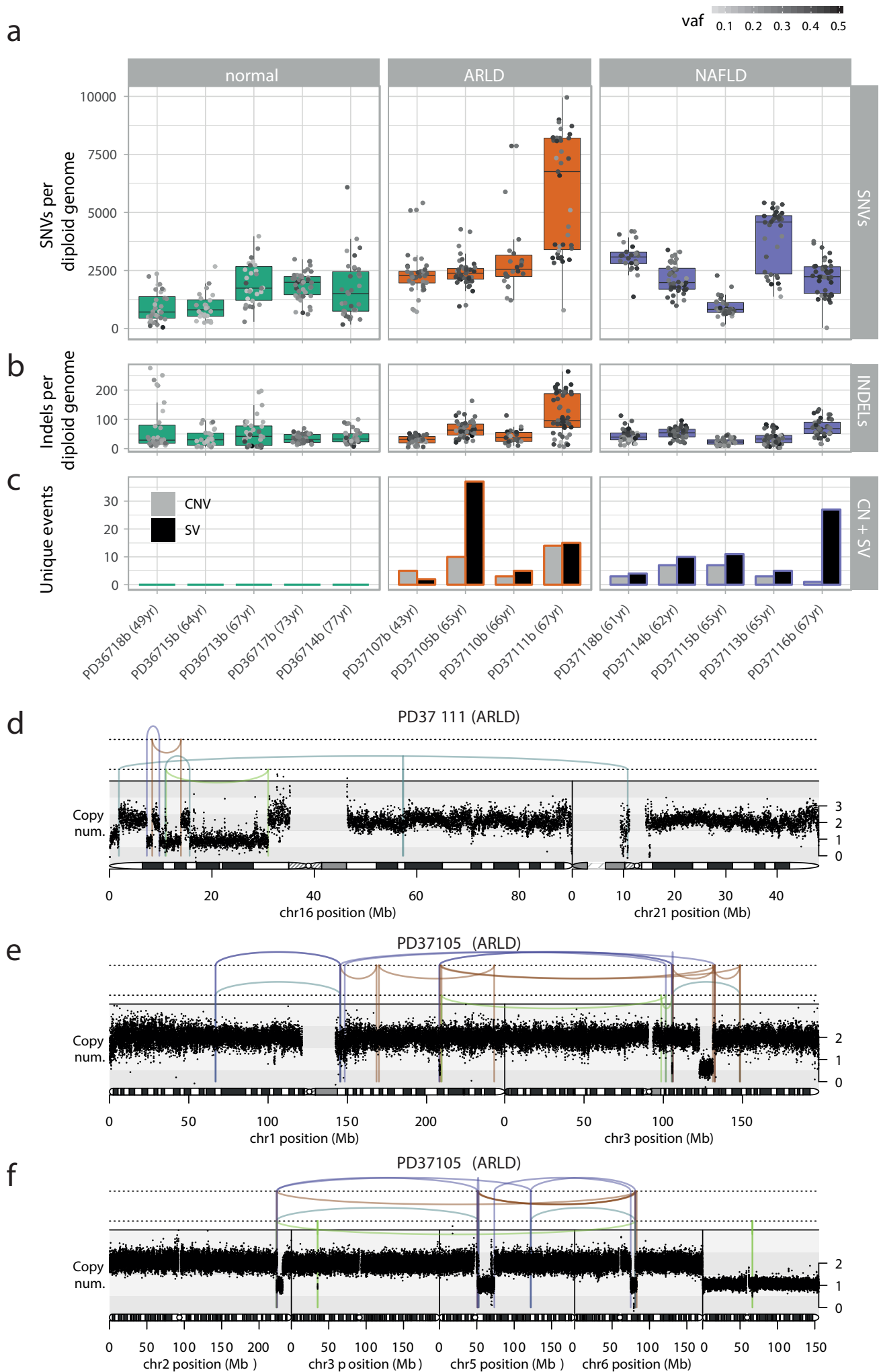


Figure 2

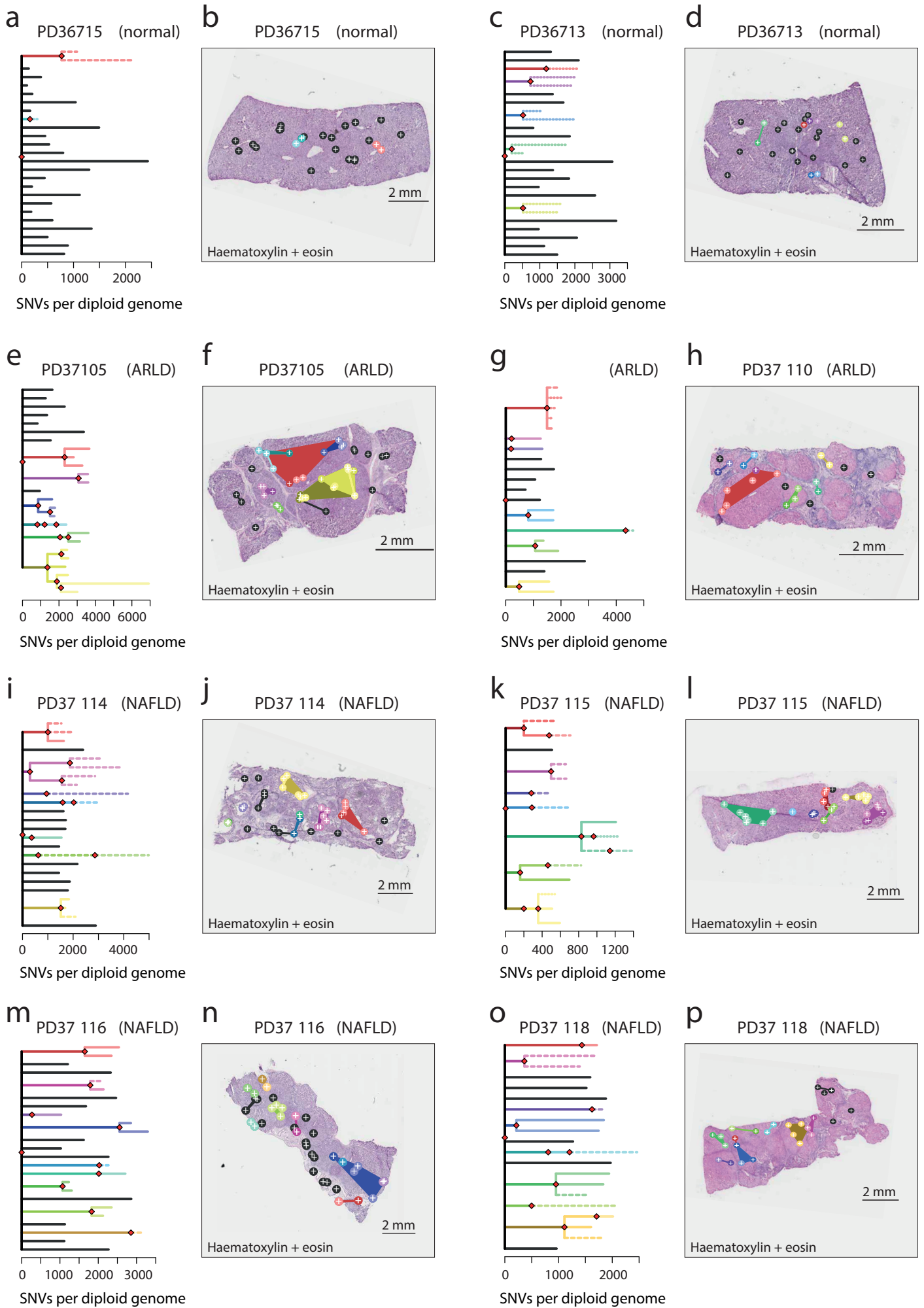


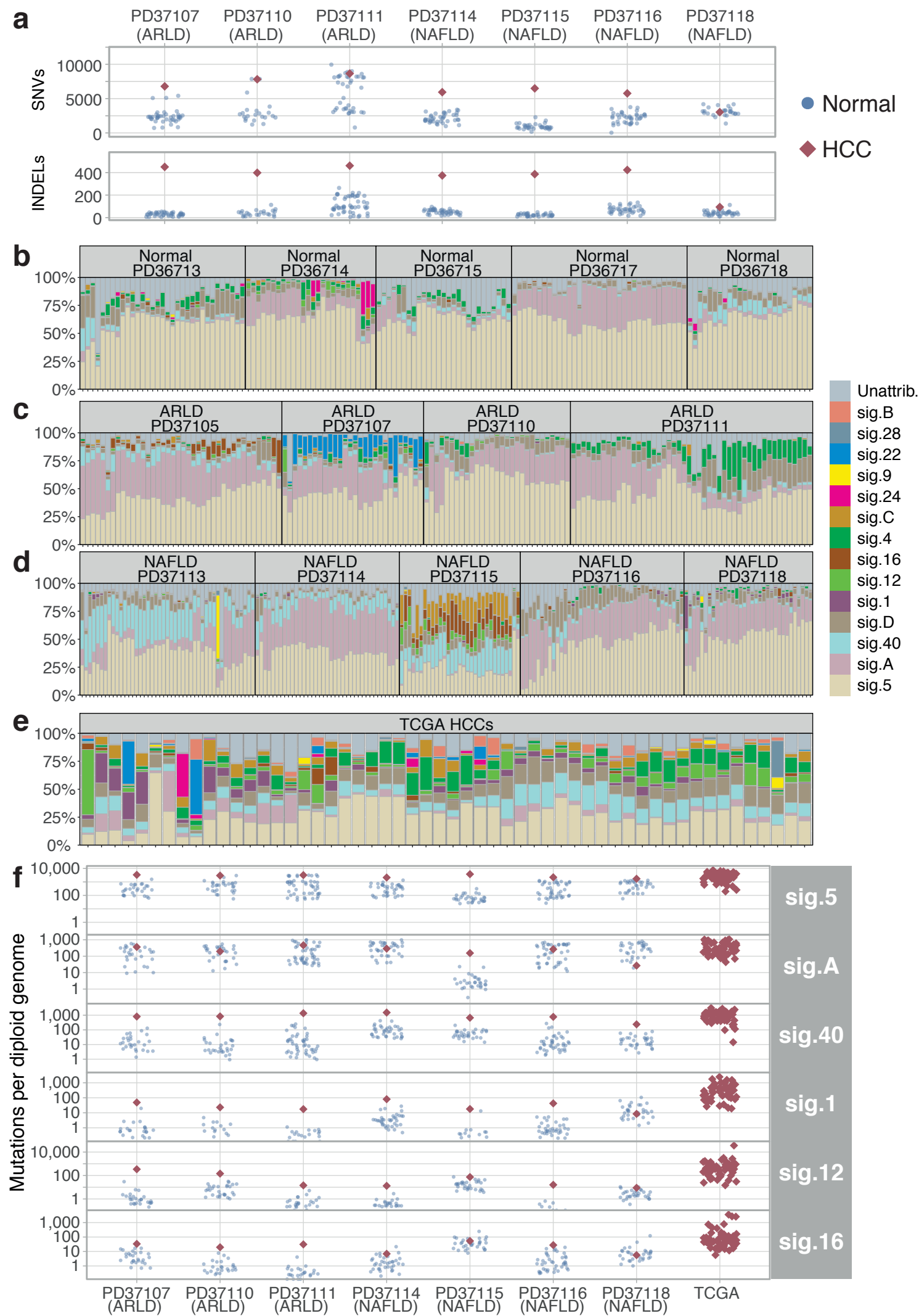
Figure 3

Figure 4

