1 **Signatures of copy number alterations in human cancer**

2 Christopher D. Steele[1], Ammal Abbasi[2,3,4], S. M. Ashiqul Islam[2,3,4], Azhar
3 Khandekar[2,3,4], Kerstin Haase[5], Shadi Hames[1], Maxime Tarabichi[5,] Tom Lesluyes[5],
4 Adrienne M. Flanagan[1,6,] Fredrik Mertens[7,8], Peter Van Loo[5], Ludmil B.
5 Alexandrov[2,3,4,*], and Nischalan Pillay[1,6,*]
6
7 [1]Research Department of Pathology, Cancer Institute, University College London,
8 London, WC1E 6BT, UK
9 [2]Department of Cellular and Molecular Medicine, UC San Diego, La Jolla, CA,
10 92093, USA
11 [3]Department of Bioengineering, UC San Diego, La Jolla, CA, 92093, USA
12 [4]Moores Cancer Center, UC San Diego, La Jolla, CA, 92037, USA
13 [5]Cancer Genomics Laboratory, The Francis Crick Institute, London, NW1 1AT, UK
14 [6]Department of Cellular and Molecular Pathology, Royal National Orthopaedic
15 Hospital NHS Trust, Stanmore, HA7 4LP, UK
16 [7]Division of Clinical Genetics, Department of Laboratory Medicine, Lund University,
17 Lund, Sweden
18 [8]Department of Clinical Genetics and Pathology, Division of Laboratory Medicine,
19 Lund, Sweden
20
21 [*]Denotes equal contributions. Correspondence and request for materials should be
22 addressed to L2alexandrov@health.ucsd.edu and N.pillay@ucl.ac.uk.
23

24 **ABSTRACT**

25 The gains and losses of DNA that emerge as a consequence of mitotic errors and

26 chromosomal instability are prevalent in cancer. These copy number alterations

27 contribute to cancer initiaition, progression and therapeutic resistance. Here, we

28 present a conceptual framework for examining the patterns of copy number

29 alterations in human cancer using whole-genome sequencing, whole-exome

30 sequencing, and SNP6 microarray data making it widely applicable to diverse

31 datasets. Deploying this framework to 9,873 cancers representing 33 human cancer

32 types from the TCGA project revealed a set of 19 copy number signatures that

33 explain the copy number patterns of 93% of TCGA samples. 15 copy number

34 signatures were attributed to biological processes of whole-genome doubling,

35 aneuploidy, loss of heterozygosity, homologous recombination deficiency, and

36 chromothripsis.  The aetiology of four copy number signatures are unexplained and

37 some cancer types have unique patterns of amplicon signatures associated with

38 extrachromosomal DNA, disease-specific survival, and gains of proto-oncogenes

39 such as *MDM2*. In contrast to base-scale mutational signatures, no copy number

40 signature associated with known cancer risk factors. The results provide a

41 foundation for exploring patterns of copy number changes in cancer genomes and

42 synthesise the global landscape of copy number alterations in human cancer by

43 revealing a diversity of mutational processes giving rise to copy number changes.

44

45 **MAIN**

46 Genome instability is a hallmark of cancer leading to changes of the genomic DNA

47 sequence, aneuploidy, and focal copy number alterations[1]. Both aneuploidy and sub-

48 chromosomal copy number alterations have been previously associated with

49 increased cell proliferation, poor prognosis, and reduced infiltration of immune cells[2–

50 6]. Aneuploidy and genome-wide structural variation may originate from mitotic

51 slippage, spindle multipolarity, and breakage-fusion-bridge (BFB) cycles[7]. Besides

52 chromosome mis-segregation, other macroevolutionary mechanisms lead to

53 changes in genomic copy number, including whole-genome doubling (WGD), where

54 the entire chromosomal content of a cell is duplicated[8] and chromothripsis where a

55 "genomic catastrophe" leads to clustered rearrangements and oscillating copy

56 number[9]. These evolutionary events may occur multiple times at different intensities

57 during tumour development leading to a highly complex genome[10–12].

58

59 The complex structural profiles of human cancers are mirrored by the intricate

60 patterns of somatic mutations imprinted on cancer genomes at a single nucleotide

61 level. Previously, we developed a computational framework that allows separating

62 these intricate patterns of somatic mutations into individual mutational signatures of

63 single base substitutions (SBS), doublet base substitutions (DBS), and small

64 insertion or deletions (ID)[13,14]. Analyses of mutational signatures have provided

65 unprecedented insights into the exogenous and endogenous processes moulding

66 cancer genomes at a single nucleotide level with mutational signatures attributed to

67 exposures to environmental mutagens, failure of DNA repair, infidelity/deficiency of

68 polymerases, iatrogenic events, and many others[15–22].

69

3

70    We recently developed a "mechanism-agnostic" approach for summarising allele-

71    specific copy number patterns in whole genome sequenced sarcomas[23] which we

72    term copy number signatures. Other cancer subtype-specific methods for

73    interrogating copy number patterns have been created and applied to ovarian cancer

74    and breast cancer[24,25]. While these initial approaches have led to biological and

75    clinical insights, there is currently no approach that allows interrogating copy number

76    signatures across multiple cancer types and across different experimental assays.

77    To address this gap we developed a new framework for deciphering copy number

78    signatures across cancer types and demonstrate its applicability to whole-genome

79    sequencing, whole-exome sequencing, and SNP6 microarray data. We identified 19

80    distinct copy number signatures many of which are shared across multiple

81    histologies and others that are specific to certain cancer subtypes.  Extensive

82    computational simulations, refinement and statistical association analyses were used

83    both to assign processes to many of these signatures and to demonstrate their

84    biological and clinical relevance. Overall, our findings shed light on the processes of

85    chromosomal segregation errors and provide a method to distil the ensuant complex

86    genomic configurations.

87

88    **A framework for pan-cancer classification of copy number alterations**

89    We examined the allele-specific copy number profiles of 9,873 primary cancer

90    samples across 33 cancer types from The Cancer Genome Atlas project (TCGA;

91    **Supplementary Table 1**). The severity of genomic instability, measured by number

92    of copy number segments, proportion of the genome displaying loss of

93    heterozygosity (LOH) and genome doubling status vary greatly amongst cancer

94    types (**Fig. 1a-b**). Nevertheless, a linear relationship was observed between the

95    number of segments and proportion of genomic LOH, varying from cancers with

96    diploid and copy number "quiet" genomes (*e.g.,* acute myeloid leukaemia, thymoma,

97    and thyroid carcinoma; **Fig. 1*a***) to cancers with highly aberrant copy number profiles

98    (*e.g.,* ovarian carcinomas and sarcomas; **Supplementary Fig. 1*a-b*).** This linear

99    relationship fails to hold only for adrenocortical carcinoma and chromophobe renal

100   cell carcinoma both of which demonstrate enrichment of LOH without enrichment of

101   copy number segmentation (**Supplementary Fig. 1*a-c***). Additionally, considerable

102   variability of ploidy was observed both between and within cancer types (**Fig. 1*b***,

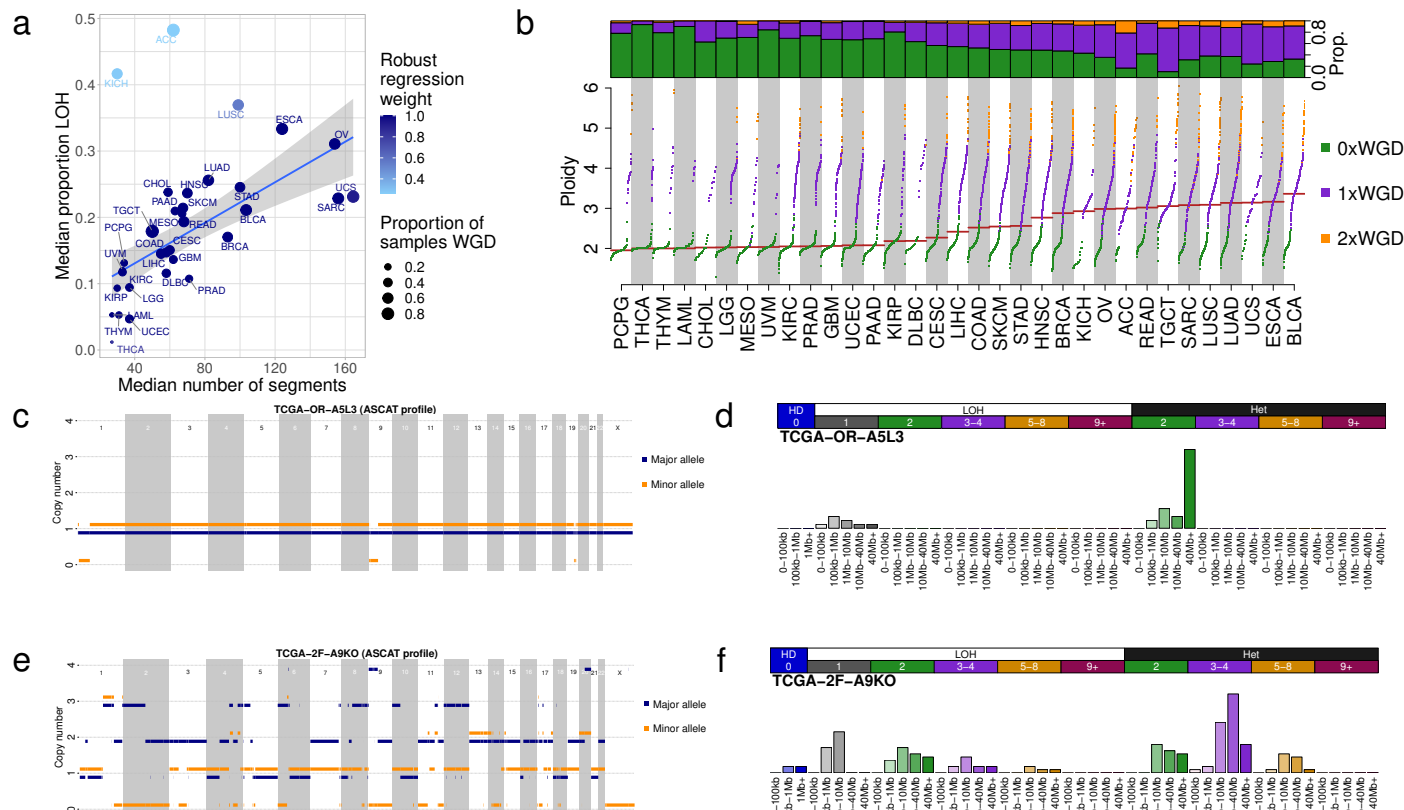103   **Supplementary Fig. 1*d***).

104

**Figure 1 – Pan-cancer copy number characteristics in TCGA.**

**a)** Copy number characteristics of 33 tumour types included in TCGA. Median number of segments in a copy number profile (x-axis), median proportion of the genome that has loss of heterozygosity (y-axis) and the proportion of samples that have undergone one or more whole genome doubling events (size). The line of best fit from a robust linear regression is shown, where the colour of points indicates the weight of the tumour type in the regression model.

**b)** Ploidy characteristics of all TCGA samples, split by tumour type. Bottom panel: ploidy (y-axis) against quantile of ploidy (y-axis) for each sample in a tumour type, where samples are coloured by their genome doubling status: 0xWGD=non genome doubled (green), 1xWGD=genome doubled (purple), 2xWGD=twice genome doubled (orange). Top panel: proportion of samples in each tumour type that are 0, 1 or 2xWGD.

**c)** Allele-specific copy number profile from a majority diploid sample (sample ID: TCGA-OR-A5L3, tumour type: ACC). Copy number (y-axis) across the genome (x-axis) is given for both the major (blue) and minor (orange) allele.

**d)** Copy number summary for TCGA-OR-A5L3 after categorizing each of the segments. Segments are characterized first as homozygously deleted (left, blue), LOH (middle, white) or heterozygous (right, black), then by copy number states: TCN=0 (blue), TCN=1 (grey), TCN=3-4 (purple), TCN=5-8 (orange) and TCN=9+ (red). Finally, segments are categorized by segment size (increasing colour saturation indicates increasing segment size): 0-100kb, 100kb-1Mb, 1Mb-10Mb, 10Mb-10Mb and 40Mb+ (bottom labels). Homozygous deletions have a largest segment size category of 1Mb+.

**e)** Allele-specific copy number profile for a highly aberrant sample (sample ID: TCGA-2F-A9KO, tumour type: BLCA).

**f)** Copy number summary for TCGA-2F-A9KO.

6

133

134    To capture biologically relevant copy number features, we developed a classification

135    framework that encodes the copy number profile of a sample by summarizing the

136    counts of segments into a 48-dimensional vector. Specifically, copy number

137    segments were classified into three heterozygosity states: heterozygous segments

138    with copy number of {$A>0$, $B>0$} (numbers reflect the counts for major allele $A$ and

139    minor allele $B$), segments with LOH with copy number of {$A>0$, $B=0$}, and segments

140    with homozygous deletions {$A=0$, $B=0$}. Segments were further subclassified into 5

141    classes based on the sum of major and minor allele (total copy number, TCN;

142    **Supplementary Fig. 1e**) and chosen for biological relevance: TCN=0 (homozygous

143    deletion), TCN=1 (deletion leading to LOH), TCN=2 (wild type, including copy-neutral

144    LOH), TCN=3 or 4 (minor gain), TCN=5 to 8 (moderate gain), and TCN>=9 (high-

145    level amplification). Each of the heterozygous and LOH total copy numbers were

146    then subclassified into five classes based on the size of their segments: 0 – 100kb,

147    100kb – 1Mb, 1Mb – 10Mb, 10Mb – 40Mb, and >40Mb (the largest category for

148    homozygous deletions was restricted to >1Mb) in order to capture focal, large scale,

149    and chromosomal copy number changes. The segment sizes were selected to

150    ensure that a sufficient proportion of segments were classified in each category

151    resulting in a reasonable representation across the pan-cancer TCGA dataset

152    (**Supplementary Fig. 1f-h**). Two examples, one encoding a mostly diploid

153    adrenocortical carcinoma (**Fig. 1c-d**) and another encoding a genomically unstable

154    bladder cancer (**Fig. 1e-f**), are provided to illustrate the classification framework.

155

156    To determine the generalizability of our framework for pan-cancer classification of

157    copy number alterations across experimental platforms, we performed a comparative

158   analysis of samples simultaneously profiled with SNP6 microarrays, whole-exome

159   sequencing (282 samples), and whole-genome sequencing (512 samples).

160   Optimisation of the copy number calling strategy (**Methods**) resulted in remarkably

161   similar profiles between distinct experimental assays. Specifically, copy number

162   profiles derived from exome sequencing data had a median cosine similarity of 0.925

163   with copy number profiles derived from SNP6 microarrays (**Supplementary Fig. 1*i***).

164   Copy number profiles derived from whole-genome sequencing data exhibited

165   median cosine similarities of 0.933 and 0.852 with profiles derived from SNP6

166   microarrays or exome sequencing, respectively (**Supplementary Fig. 1*j-k***). These

167   similarities are considerably better than similar comparisons observed for mutational

168   signatures of single base substitutions derived from whole-genome and exome

169   sequencing (median cosine similarity=0.55).

170

171   **The repertoire of copy number signatures in human cancer**

172   Copy number profiles from SNP6 microarrays (n=9,873) were concatenated into

173   cancer type-specific matrices and separately in a global pan-cancer matrix. These

174   matrices were decomposed using our previously established approach[26] for deriving

175   a reference set of signatures (**Methods**). The approach allowed the identification of

176   both the shared patterns of copy number across all examined samples, termed, *copy*

177   *number signatures*, as well as the quantification of the number of segments

178   attributed to each copy number signature in each sample, termed, *signature*

179   *attribution*.

180

181   By applying our copy number signature framework (**Methods**) we identified 19

182   distinct pan-cancer signatures (**Fig. 2*a***; **Supplementary Table 2**). These signatures

8

183   accurately explained the copy number profiles (p-value<0.05, Methods) of 93% of

184   the examined TCGA samples. The remaining 7% were poorly explained due to a

185   combination of a low number of segments and/or a high diversity of copy number

186   states in the copy number profile or few operative signatures identified

187   (**Supplementary Figs. 2*a-c***). The 19 signatures were categorized into 6 groups

188   based on their most prevalent features. CN1 and CN2 are primarily defined by

189   >40Mb heterozygous segments with total copy number (TCN) of 2 and 3-4

190   respectively. CN3 is characterized by heterozygous segments with sizes above 1Mb

191   and TCN between 5 and 8. CN4-8 each have segment sizes between 100kb and

192   10Mb but with different TCN or LOH states. CN9-12 each have numerous LOH

193   components with segment size <40Mb. CN13-14 have whole-arm or whole-

194   chromosome scale LOH events (>40Mb). CN15 consists of LOH segments with TCN

195   between 2 and 4 as well as heterozygous segments with TCN between 3 and 8,

196   each with segment sizes 1-40Mb. CN16-19 exhibited complex patterns of copy

197   number alterations that are uncommon but are seen in distinct cancer types.

198   Additionally, 3 artefactual signatures (CN20-22) indicative of copy number profile

199   over-segmentation were identified (**Supplementary Fig. 2*d***). To determine if the

200   copy number signatures would generalize between platforms, we compared copy

201   number signatures derived from whole-genome and whole-exome sequencing with

202   SNP6 array signatures which showed a strong concordance with a median cosine

203   similarity between signatures above 0.80 (**Supplementary Fig. 2*e-h***).
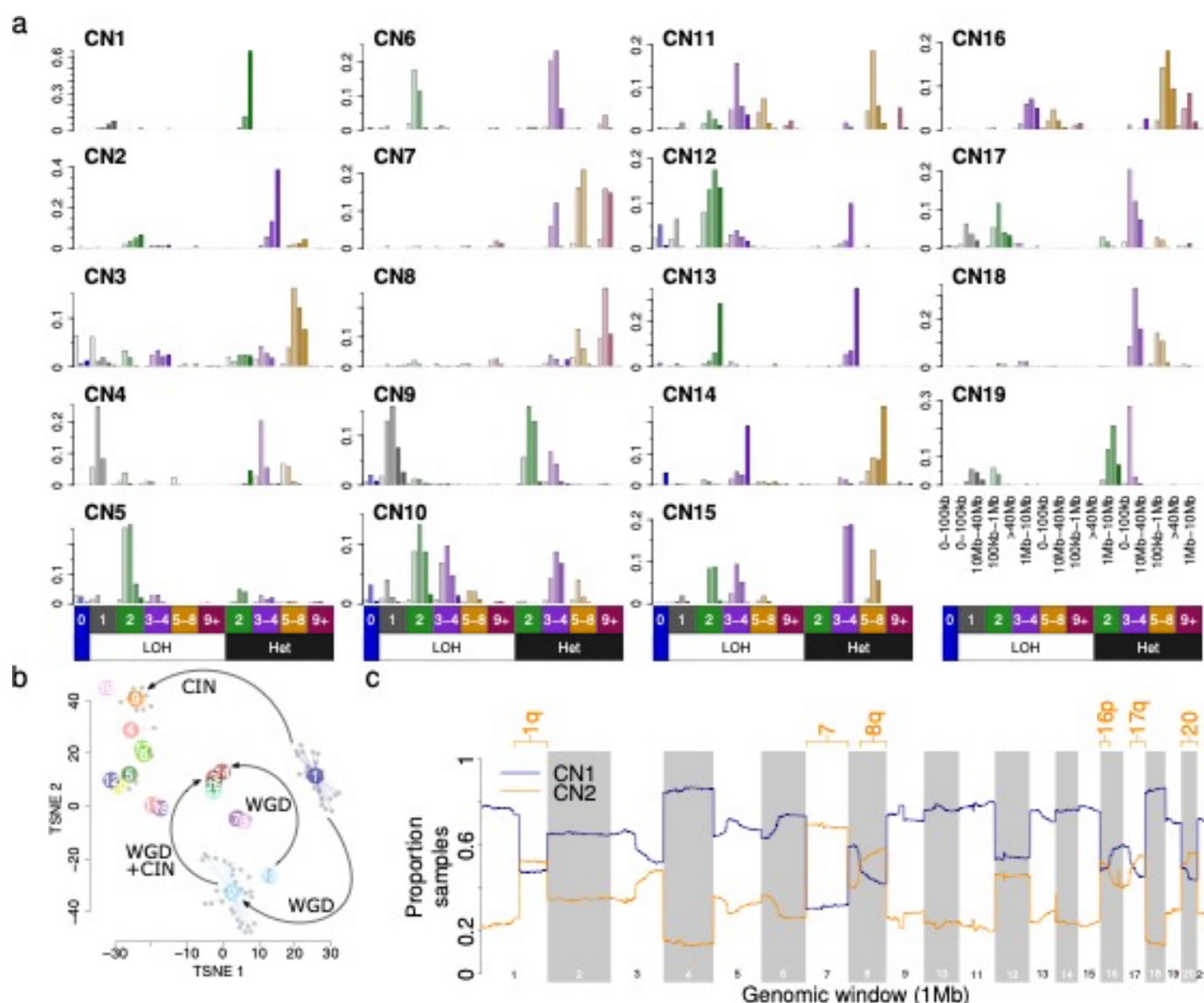
204

**Figure 2 – Patterns of pan-cancer copy number signatures.**

a) 19 identified non-artefactual copy number signatures in TCGA that are not linear combinations of any other. LOH status and total copy number are indicated below each column. Segment sizes for select bars are shown in the bottom right. Increasing saturation of colour indicates increasing segment size.

b) TSNE representation of all non-artefactual consensus signatures (colours) and the individual signatures that were combined to form each consensus signature (grey). Inferences about the relationships between signatures (see Supplementary Figure 3) are indicated with arrows; WGD=whole-genome doubling, CIN=chromosomal instability.

c) CN1 (blue) and CN2 (orange) recurrence (y-axis) across the genome (x-axis) in 472 highly aneuploid samples where CN1+CN2 attribution = 1. Chromosome arms with >50% samples attributed to CN2 are labelled.

10

220

**The transitional behaviour of copy number signatures**

222 The catalogue of somatic mutations of a cancer genome is the cumulative result of

223 the mutational processes that have been operative over the lifetime of the cell from

224 which the cancer has derived[27]. Analysis of SBS and ID mutational signatures have

225 used assumptions and prior evidence that individual mutations are independent and

226 additive[28]. However, this assumption is clearly violated for large-scale macro-

227 evolutionary events such as whole-genome doubling[29].

228

229 We therefore generated several synergistic lines of evidence to investigate the

230 impact of genome doubling on copy number signatures. First, each copy number

231 signature was tested for enrichment in non-, once- or twice-genome doubled

232 samples (**Supplementary Fig. 3a-b**). Second, *in silico* simulations of genome

233 doubling on the extracted signatures were performed (**Methods**; **Supplementary**

234 **Fig. 3c**). Third, copy number profiles arising from dynamics of whole-genome

235 doubling and chromosomal instability (CIN) were simulated (**Supplementary Fig.**

236 **3d**) and re-examined for the previously derived signatures (**Supplementary Fig. 3e**).

237

238 By combining the preceding set of experiments, we revealed a transitional behaviour

239 of copy number signatures with one signature being completely replaced by another

240 upon genome doubling (**Fig. 2b**). In this model, a cancer with a diploid signature

241 (CN1), may undergo genome doubling, thus altering signature CN1 into signature

242 CN2, or may undergo chromosomal instability transforming signature CN1 into

243 signature CN9. Through a combination of CIN and genome doubling CN2 may also

11

244     be changed to CN3. Additionally, CN13 and CN14 may be linked through genome

245     doubling, on the background of early chromosomal losses.

246

247     While macro-evolutionary events have a transitional effect on copy number

248     signatures, we hypothesized that smaller-scale events, such as segmental

249     aneuploidy, may reflect an additive behaviour. To investigate this, we focused on the

250     ploidy-associated signatures CN1 and CN2, where a combination of both signatures

251     indicates a hyper-diploid or sub-tetraploid profile. Interestingly, each signature was

252     found at below 50% attribution in approximately a quarter of TCGA samples,

253     suggestive of potential aneuploidy in a considerable proportion of samples. We

254     mapped these signatures across the cancer genomes with mixtures of attributions

255     from signatures CN1 and CN2 (**Supplementary Fig. 3*f***). This analysis recapitulated

256     known patterns of aneuploidy in human cancer[30,31], including gains of chromosomes

257     1q, 7, 8q, 16p, 17q, and 20 in more than 50% of TCGA samples (**Fig. 2*c***).

258

259     **The landscape of copy number signatures**

260     Next, we surveyed the distribution of the 19 signatures across the different cancer

261     types (**Fig. 3**). Unsurprisingly, the ploidy associated signatures CN1 and CN2 were

262     found in most samples across all cancer types with different median attributions.

263     Signatures CN4, CN7, CN10, CN16, CN18, and CN19 were derived through cancer

264     type extractions and therefore unique to uveal melanoma, breast cancer, lung

265     squamous carcinoma, ovarian carcinoma, liver cancer and paragangliomas,

266     respectively. Signatures CN4-8 all showed segments of high total copy number and

267     were seen in tumour types with known prevalent amplicon events[32]. CN9-CN12

268     showed differing patterns of hypodiploidy, LOH < 40Mb and WGD reflective of

269    chromosomal instability. Signatures CN13 and CN14 were prevalent in

270    adrenocortical carcinoma and chromophobe renal cell carcinoma, suggesting a link

271    with the known patterns of chromosomal LOH (cLOH) seen in these cancers[33,34].

272    Signature CN15 was prevalent in tumour types previously described as being

273    enriched in the tandem duplicator phenotype (TDP)[35]. Different cancer lineages

274    clustered together based on the prevalence of signatures; namely TDP, whole-

275    genome duplication, diploid chromosomal instability, simple diploidy, and

276    chromosomal LOH (**Fig. 3**). This segregation of cancer types and their constituent

277    signatures reflects the known distributions of genome doubling and aneuploidy in
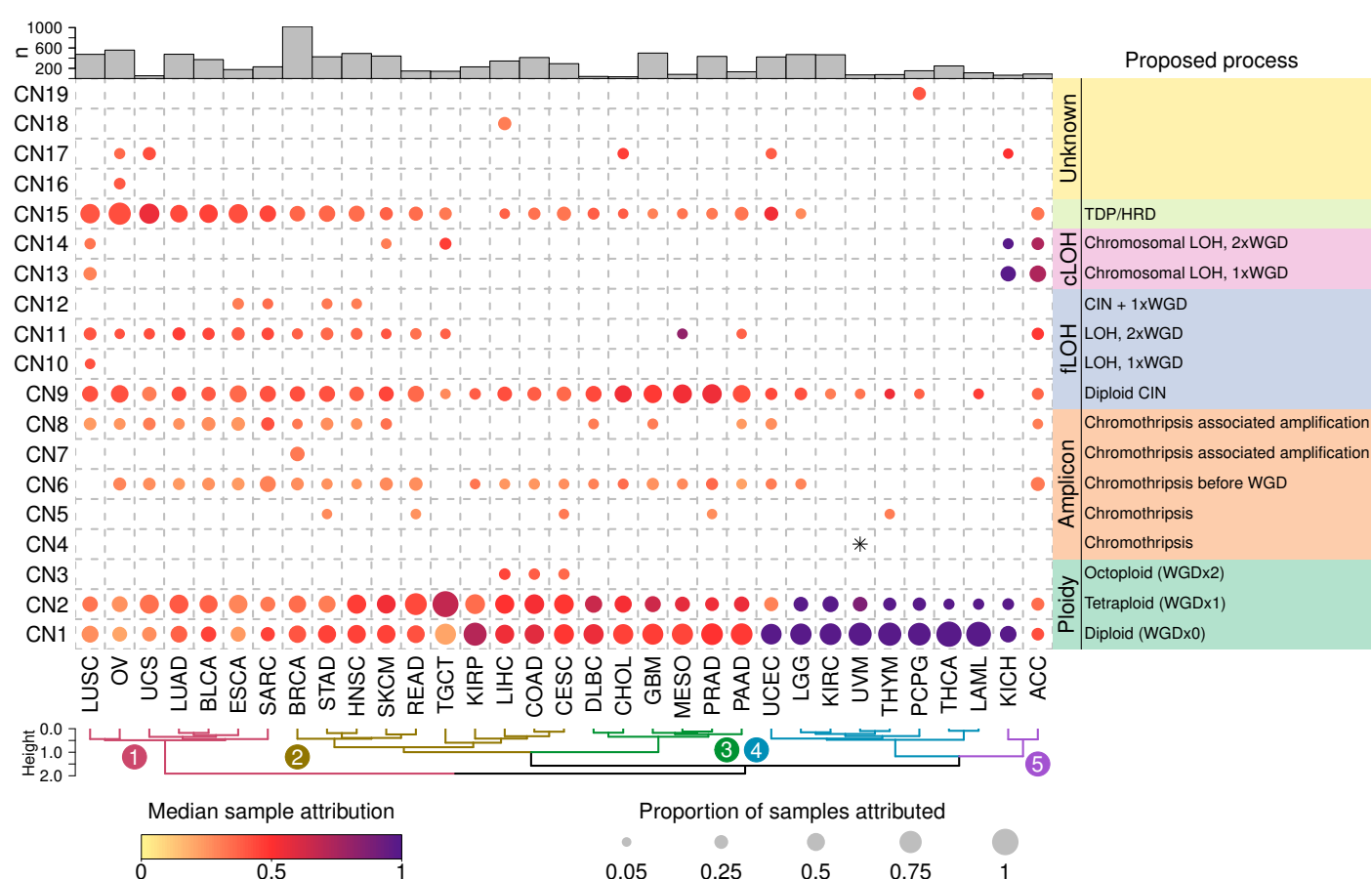
278    human cancer[3,36].

279



281    **Figure 3 – Distribution of copy-number signatures across human cancers.**
282    Attributions of the 19 non-artefactual signatures (y-axis) split by tumour type (x-
283    axis), showing both the proportion of each tumour type exposed to each

13

284        signature (size), and the median exposure of those tumours that are exposed
285        to the signatures in each tumour type (colour). Tumour/signature combinations
286        with less than 5% of samples exposed to the signature are not shown (except
287        for CN4 in UVM, denoted with a *). Hierarchical clustering is shown below,
288        sample sizes are shown above. Proposed processes are shown to the right.

289

290    **Copy number signatures associated with amplicons**

291    Oncogene amplification has been associated with aggressive behaviour in cancer[32],

292    and can originate through the processes of BFB cycles and chromothripsis[12,37].

293    Reasoning that signatures with high levels of total copy number (CN4, CN5, CN6,

294    CN7, and CN8) could associate with genomic amplification we correlated these

295    signatures with known classes of amplicons[32,38]. All amplicon signatures were

296    positively associated with one or more amplicon types (**Fig. 4a**); CN8 was strongly

297    associated with all four classes of amplicon, but most strongly with extra-

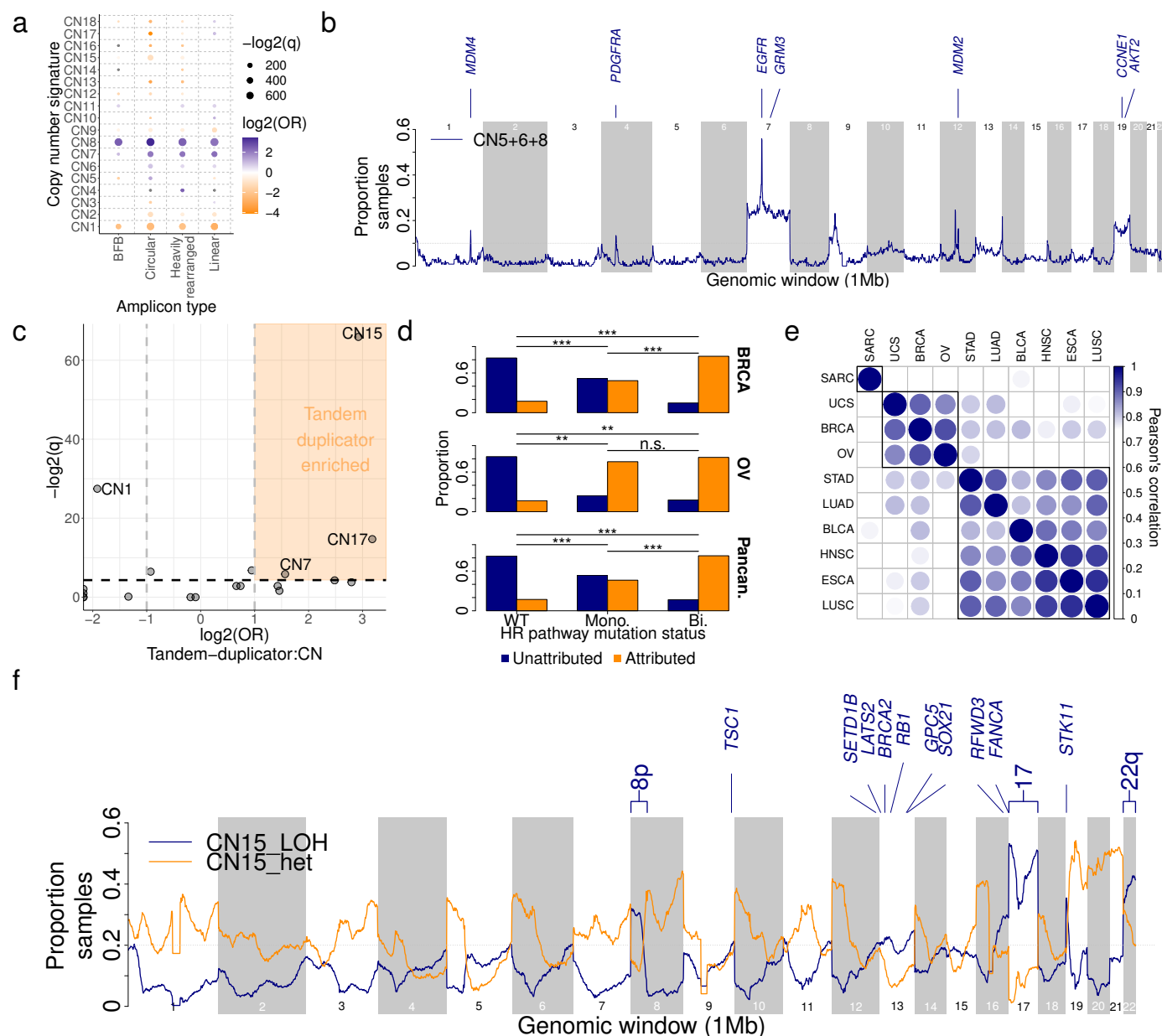298    chromosomal circular DNA amplicons (ecDNA).

299

14

**Figure 4 – Biological inference of copy-number signatures.**

**a)** Associations between copy number signatures (y-axis) and amplicon structures (x-axis), displaying the q-value (size) and log2 odds ratio (colour) from a Fisher's exact test of genomic regions attributed/not attributed to each signature against each amplicon type. Non-significant (q≥0.05) associations are not shown. BFB=breakage fusion bridge. CN8 was most strongly associated with circular amplicons: OR=10.8, q<5e-324.

**b)** Recurrence of mapped amplicon signatures (CN5, CN6 and CN8) in 1Mb windows of the human genome across 134 GBM in which the amplicon signatures were attributed. Oncogenes in regions with >10% samples attributed to amplicon signatures are labelled.

**c)** Associations between copy number signature attributed samples and tandem-duplicator phenotype samples, displaying -log2(q-values) (y-axis) and log2 odds ratios (x-axis). CN15 association: OR=7.6, q=1.5e-20, Fisher's exact test.

**d)** Correlation of CN15 attribution (y-axis) with mutational status of one or more genes of the homologous recombination pathway (x-axis) in breast cancer

15

317         (top), ovarian cancer (middle) or pan-cancer (bottom). WT=wild type. Mono =

318         Mono-allelic and Bi = bi-allelic. *=q<0.05, **=q<0.01, ***=q<0.001, n.s.=q≥0.05.

319     **e)**   Pearson's correlation of recurrence of mapping of LOH segments of CN15 to

320         the genome calculated for all pairwise comparisons of CN15-enriched tumour

321         types.

322     **f)**   Recurrence of mapped CN15 in 1Mb windows of the human genome in all

323         CN15 attributed BRCA, OV and UCS samples, split by LOH (blue) and

324         heterozygous segments (orange). Tumour-suppressor genes in regions with

325         >20% samples attributed to CN15 with LOH segments are labelled.

326

327   Recent evidence revealed that genomic amplification can evolve through interrelated

328   processes of chromothripsis, BFB and ecDNA formation[11]. Therefore, we mapped

329   the CN signatures with known regions of chromothripsis[39] across the genome

330   (**Methods**), revealing CN5-8 as being enriched in chromothriptic regions

331   (**Supplementary Fig. 4a**). Each of these signatures are dominated by small

332   segments, while CN7-8 are both strongly associated with amplified chromothripsis[40]

333   (**Supplementary Fig. 4b**) and complex chromothriptic events (**Supplementary Fig.**

334   **4c**). Simulations of copy number profiles incorporating processes of chromothripsis,

335   whole-genome doubling, and chromosomal duplication (**Supplementary Fig. 4d**)

336   demonstrated that CN4 to CN8 can be generated through chromothripsis-like events,

337   and that these signatures reflect distinct life histories of tumours, such as

338   chromothripsis before or after genome doubling (**Supplementary Figs. 4c & e**).

339

340   Chromothripsis and gene amplification are both independently associated with poor

341   prognosis[32,41]. Attribution of any of the five amplicon signatures in their respective

342   cancer types resulted in a poor disease-specific survival in a univariate pan-cancer

343   analysis (**Supplementary Fig. 5a**). Similarly, multiple amplicon signatures were

344   associated with a reduced disease-specific survival in multivariate pan-cancer and

345   cancer type analyses with consistent results from analyses based on Cox-model

346   hazard ratios (**Supplementary Fig. 5b-c**) and analyses based on accelerated failure

347    times (**Supplementary Fig. 5d-e**). Cancer type-specific survival analysis revealed

348    that patients with glioblastoma with operative signature CN5 had a poor disease-

349    specific survival (172 days reduced median survival; **Supplementary Figure 5d**). To

350    determine the topographic localization of the amplification events, we mapped the

351    amplicon signatures operative in glioblastoma (CN5, CN6, and CN8) across the

352    genome which revealed recurrence of regions involving *EGFR, PDGFRA* and *MDM2*

353    (**Fig. 4b**) in keeping with previous reports of chromothripsis-associated amplification

354    of these genes[42].

355

356    **Copy number signatures associated with loss of heterozygosity**

357    Loss of heterozygosity (LOH) is an important mechanism contributing to the

358    inactivation of tumour suppressor genes during cancer development[39,43,44]. We found

359    that 7 signatures positively correlated with LOH regions of the genome

360    (**Supplementary Fig. 6a**). Four of these signatures (CN9-12) were designated focal

361    LOH (fLOH) signatures as they exhibited predominant segments sizes <40Mb (**Fig.**

362    **2**). The four fLOH signatures were recurrently found around tumour suppressor

363    genes (**Supplementary Fig. 6b**).

364

365    In adrenocortical carcinoma and chromophobe renal cell carcinoma a characteristic

366    pattern of chromosome-level LOH leads to hypodiploidy[45,46]. We identified 2

367    signatures (CN13 and CN14) of chromosomal-scale LOH, each of which was

368    enriched in both of these cancers (**Supplementary Fig. 6c-d**). Mapping of these

369    signatures to the genome revealed recurrent LOH in chromosome regions 1p, 3p,

370    5q, 9, 10q, 13q, and 17p (**Supplementary Fig. 6e**), matching known patterns of

371    aneuploidy in these tumours[33,34] (**Supplementary Fig. 6f-g**).

372

**Copy number signature associated with tandem duplication and homologous**

**recombination deficiency**

375 Somatic tandem duplications (TD) are commonly found in breast and ovarian

376 cancer[35,47,48]. Further, TD are strongly associated with failure of homologous

377 recombination repair of DNA double strand breaks e.g. due to defective *BRCA1* or

378 *BRCA2*[35,47,48]. A detailed characterization of TD across cancer has revealed three

379 patterns with duplicated segments[35] ranging around 10kb, 200kb, or 2Mb,

380 respectively. CN15 has a segment size distribution that overlaps with the largest of

381 these three patterns and was strongly associated with TD (**Fig. 4c**, OR=7.6, q=1.5e-

382 20, Fisher's exact test) and enriched in cancer types known to show TD

383 (**Supplementary Fig. 7a**)[35].

384

385 Consistent with prior observations for TD, an enrichment of CN15 is observed for

386 samples harbouring mono-allelic defects in the homologous recombination pathway

387 compared to wild-type samples for breast cancer (**Fig. 4d**; OR=4.5 with q=6.1e-14;

388 Fisher's exact test), ovarian cancer (OR=15.3 with q=5.9e-3), and across all cancers

389 (OR=4.2 with q=2.2e-106). Further enrichments of CN15 were observed in samples

390 with bi-allelic defects in the homologous recombination pathway compared to

391 samples with mono-allelic defects for breast cancer (**Fig. 4d**; OR=6.2 with q=6.2e-5;

392 Fisher's exact test) and across all cancers (OR=5.7 with q=4.3e-16).

393

394 Prior analysis has shown that breakpoints resulting from TDs segregate non-

395 randomly in the genome[35]. Mapping of CN15 to the genomes of CN15-enriched

396 cancers revealed a tumour type-specific distribution of LOH segments (**Fig. 4e**), but

397    not of heterozygous segments (**Supplementary Fig. 7b**). Breast and ovarian cancer

398    as well as uterine carcinosarcoma displayed recurrent chromosomal LOH at 8p, 17

399    (including *BRCA1* and *TP53)*, and 22 (**Fig. 4f**). Focal LOH was also observed on 9q

400    around *TSC1*, 13q around *BRCA2* and *RB1*, and 19p around *STK11* (**Fig. 4f**). In

401    contrast CN15 attributed sarcomas display strong peaks of recurrent LOH around

402    known sarcoma tumour suppressor genes[49] (*CDKN2A*, *RB1*, and *TP53*;

403    **Supplementary Fig. 7c**). The 6 other tumour types enriched in CN15 display

404    recurrent chromosomal LOH at 8p, 9p, 17p, 19p, and 21 (**Supplementary Fig. 7d**).

405

406    **Copy number signatures associate with genomic features**

407    To identify DNA damage repair mechanisms involved in the mutational processes

408    giving rise to copy number signatures, we evaluated the associations between the

409    activities of copy number signatures  and single nucleotide level mutational

410    signatures from both exome and whole genome sequencing data (**Fig. 5a).** As

411    previously described SBS3 and ID6 are strongly associated with defective

412    homologous recombination repair[14]. SBS2 and SBS13 are associated with

413    APOBEC-mediated mutagenesis particularly seen near double stranded DNA

414    breaks[50]. As expected, CN15 was strongly associated with SBS3 and ID6 derived

415    from both WES and WGS data. Additionally, CN15 was associated with SBS2 and

416    SBS13 providing a putative mechanistic link between APOBEC activity and CN15 in

417    the context of TDPs. Negative associations were observed for diploid signature CN1

418    and APOBEC signatures SBS2 and SBS13 as well as for CN1 and tobacco-

419    associated signature SBS4. These results indicate that diploid cancer genomes have

420    lower APOBEC mutagenesis and that most cancers of tobacco smokers are not

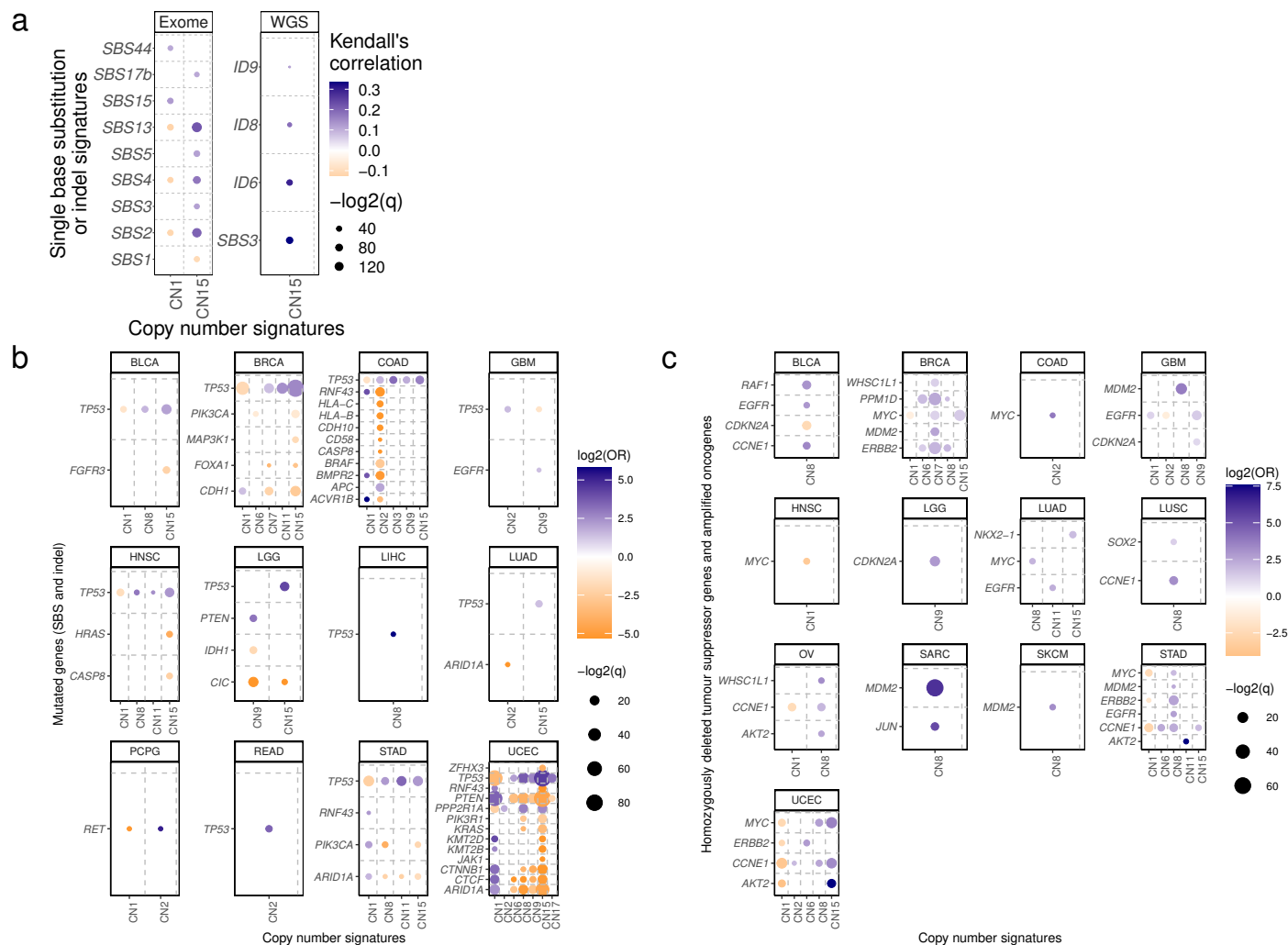421    diploid.

19

422

**Figure 5 – Genomic associations of copy number signatures.**

**a.** Correlation between copy number signature (y-axis) attribution and single base substitution signature (x-axis, SBS) exposure across TCGA exomes (left) and whole genomes (right). Strength of correlation is indicated by colour (orange=anti-correlated, blue=correlated), q-value is indicated by size of point. Only SBS signatures with any correlation between any copy number signatures with q<0.01 are shown. CN15 association with exome SBS3: Kendall's correlation=0.12, q=7.5e-12. CN15 association with exome SBS2 and SBS13: Kendall's correlation=0.2 and 0.22, q=1.6e-43 and 2.2e-50, respectively. CN15 association with WGS SBS3: Kendall's correlation=0.34, q=1.1e-21. CN15 association with WGS ID6: Kendall's correlation=0.29, q=4.7e-15.

**b.** Associations between copy number signatures (x-axis) and driver gene SNV/indel status (y-axis) across each TCGA tumour type (panels). Effect size (log2 odds ratio, colour), and significance level (-log2 q-value, size) from a Fisher's exact test are displayed.

**c.** Associations between copy number signatures (x-axis) and driver gene copy number alteration status (y-axis, amplification for oncogenes, homozygous deletion for tumour-suppressor genes) across each TCGA tumour type (panels). Effect size (log2 odds ratio, colour), and significance level (-log2 q-value, size) from a Fisher's exact test are displayed.

21

445   We next interrogated cancer driver gene mutations and copy number signatures and

446   found significant differences between cancer types. A consistent finding across

447   cancer was a positive association between *TP53* mutation and multiple copy number

448   signatures (**Fig. 5b**). *TP53* mutations were also associated with an increased

449   diversity of copy number signatures (**Supplementary Fig. 8a**; OR=3.42 with q=1.5e-

450   49), supporting the link between *TP53* alteration and aneuploidy[3,51–53]. Mutations in

451   *RNF43, HLA-B, HLA-C* and *BRAF* are commonly seen in microsatellite instable

452   (MSI) colon cancers and were found to be negatively correlated with samples with

453   tetraploid genomes (*i.e.,* CN2 attributed; **Supplementary Fig. 8b**). MSI is associated

454   with high immune cell infiltration whilst aneuploidy is associated with a decrease in

455   leucocyte fraction[54]. Across multiple cancer types, we observe a general trend of

456   decreased leucocyte fractions in cancers with copy number signatures of aneuploidy

457   compared to diploid cancers (CN1; **Supplementary Fig. 8c**). Similar to colon

458   cancer, multiple cancer driver genes were associated with CN1/CN2 in endometrial

459   cancer, largely driven by differential copy number and mutation patterns seen in

460   microsatellite stable and unstable tumours (**Supplementary Fig. 8d**).

461

462   To assess the relationships between copy number signatures and copy number

463   driver genes, we evaluated the associations between attributions of copy number

464   signatures and homozygous deletions of COSMIC tumour suppressor genes as well

465   as between attributions of copy number signatures and amplifications of known

466   proto-oncogenes[55]. Copy number drivers such as *MDM2, EGFR, CCNE1, MYC*, and

467   *ERBB2* were strongly positively associated with amplicon signatures CN6-8 as well

468   as CN15 (**Fig. 5c**). In contrast, *CDKN2A* was the only homozygously deleted tumour

469   suppressor gene associated with any signature, most commonly CN9.

22

470

471    In contrast to single-nucleotide level SBS and ID signatures[14], no associations were

472    found between any copy number signature and cancer risk factors: gender, smoking

473    status, or alcohol consumption (**Supplementary Fig. 8e**). Significant associations

474    were found between age and copy number signature attribution in individual tumour

475    types (**Supplementary Fig 8f**), however, these were driven by tumour sub-type

476    differences: serous *versus* endometrioid endometrial cancers (difference in mean

477    age at diagnosis=4.7 years, p=9.0e-5, Mann-Whitney test) in which non-

478    endometrioid endometrial cancers are strongly associated with HRD[56] and enriched

479    in CN15 (OR=16.7, p<7.1e-26, Fisher's exact test); synovial sarcoma *versus* other

480    sarcoma (difference in mean age at diagnosis=-22.3 years, p=4.3e-3, Mann-Whitney

481    test) in which synovial sarcomas are karyotypically simple[49] and enriched in CN1

482    (OR=Inf, p=2.3e-5, Fisher's exact test).

483

**DISCUSSION**

484    In this report, we provide the first pan-cancer framework for analysing copy number

485    signatures as well as the first comprehensive analysis of copy number signatures in

486    human cancer. The results revealed multiple distinct copy number signatures

487    including ones attributed to ploidy, amplification, loss of heterozygosity,

488    chromothripsis, and tandem duplications. Multiple signatures of unknown processes ,

489    cancer subtype specific signatures as well as artefactual signatures were identified.

490    Unlike SBS and ID mutational signatures, copy number signatures did not associate

491    with known cancer risk factors. Rather, copy number signatures reflect the activity of

492    endogenous mutational processes such as homologous recombination deficiency,

493    aberrant mitotic DNA replication, and chromothripsis[11,12].

495

496    The field of copy number signatures is nascent, with three distinct methods

497    previously implemented in three distinct tumour types[23–25]. As the field matures it will

498    become increasingly clear which models are better suited to addressing specific

499    clinical or biological questions. To resolve these questions, pan-cancer analyses

500    utilizing all of these methods will be key, and we present here the first step towards

501    that goal; a mechanism-agnostic pan-cancer compendium of allele-specific copy

502    number signatures.

503

504    **ACKNOWLEDGEMENTS**

24

534

**CONFLICTS OF INTEREST**

LBA is an inventor of a US Patent 10,776,718 for source identification by non-

negative matrix factorization.

538

**AUTHORS CONTRIBUTIONS**

Study was conceived and designed by CDS, NP and LBA. Data analysis was

performed by CDS, AA, SMAI, AK, KH, SH, MT and TL. Manuscript was written by

CDA, NP and LBA. Interpretation of data and contributions to writeup were provided

by MT, TL, AMF, FM and PVL.

544

**DATA AVAILABILITY**

No new data was generated for this study. ASCAT copy number profiles that were

generated for a different study and analysed here can be found at:

https://github.com/VanLoo-lab/ascat/tree/master/ReleasedData/TCGA_SNP6_hg19

549

**CODE AVAILABILITY**

Code for summarising copy number profiles into 48-length vectors can be found at:

https://github.com/AlexandrovLab/SigProfilerMatrixGenerator

Code for extracting copy number signature can be found at:

https://github.com/AlexandrovLab/SigProfilerExtractor

Code for decomposing copy number summaries into known copy number signatures

can be found at:

https://github.com/AlexandrovLab/SigProfilerSingleSample

Bespoke scripts for all other analysis are available from authors upon request.

559

26

560    **ONLINE METHODS**

561    **Utilized datasets**

562    Using SNP6 microarray data, copy number profiles were generated for 9,873

563    cancers and matching germline DNA of 33 different types from The Cancer Genome

564    Atlas (TCGA)[43] using allele-specific copy number analysis of tumours (ASCAT)[58]

565    with a segmentation penalty of 70 (**Supplementary Table 1**). Additionally, a set of

566    whole-genome sequences from 512 cancers of the International Cancer Genome

567    Consortium (ICGC) that overlapped with tumour profiles in TCGA were analysed[39] to

568    generate WGS-derived copy number profiles(see below). Lastly, a set of whole-

569    exome sequences from 282 cancers from TCGA was analysed to generate exome-

570    derived copy number profiles (see below).

571

572    **Copy number profile summarization**

573    Copy number segments were categorized into three heterozygosity states:

574    heterozygous (CN={>0,>0} for the major and minor alleles respectively), loss of

575    heterozygosity (LOH; CN={>0,0}) and homozygous deletion (CN={0,0}). Segments

576    were further subclassified into 5 categories of total copy number: CN0 reflects

577    homozygous deletions, CN1 represents a genomic deletion, CN2 represents a

578    diploid state, CN3-4 is a tri-to-tetraploid or gained state, CN5-8 is a penta-to-

579    octoploid state and CN9+ represents high-level amplifications. Segments were

580    further subclassified into 5 size categories: 0-100kb, 100kb-1Mb, 1Mb-10Mb, 10Mb-

581    40Mb, and >40Mb. For homozygous deletions only 3 size categories were used: 0-

582    100kb, 100kb-1Mb, and >1Mb. In this way copy number profiles were summarized

583    as counts of 48 combined copy number categories defined by heterozygosity, copy

584    number and size, which we will define as $N = [n_1, n_2, ..., n_{48}]$. For a given dataset,

585    the copy number profiles of a set with $S$ samples are then summarized as a

586    nonnegative matrix with $S \times 48$ dimensions.

587

588    **Deciphering signatures of copy number alterations**

589    Copy number signatures were extracted by applying our previously developed

590    approach for creating a reference set of signatures[14]. Specifically,

591    SigProfilerExtractor v1.0.17[26] was applied to the matrix encompassing all TCGA

592    samples as well as separately to each matrix corresponding to an individual tumour

593    type. In brief, SigProfilerExtractor utilizes nonnegative matrix factorization (NMF) to

594    find a set of copy number signatures ranging from 1 to 25 components for each

595    examined matrix. For each number of components, 250 NMF replicates with distinct

596    initializations of the lower dimension matrices were performed on the Poisson

597    resampled data. SigProfilerExtractor was used with default parameters, except for

598    the initializations of the lower dimension matrices where random initialization was

599    utilized consistent with our prior analyses of mutational signatures[14,59] After

600    performing 250 nonnegative matrix factorizations, SigProfilerExtractor clusters the

601    factorization within each decomposition to automatically identify the optimum number

602    of operative signatures that best explain the data without overfitting these data[26].

603

604    As previously done[60], the sets of all identified copy number signatures were

605    combined into a reference set of pan-cancer copy number signatures by leveraging

606    hierarchical clustering based on the cosine dissimilarities between each signature.

607    The number of combined signatures is chosen to maximise the minimum average

608    cosine similarity between each signature in a cluster and the mean of all samples in

609    that cluster, to ensure that each copy number signature in a cluster has a high

610    similarity to the combined copy number signature for that cluster. Simultaneously,

611    the maximum cosine similarity between mean copy number signatures for each

612    cluster is minimized, to ensure that each combined signature is distinct from all

613    others. To avoid reference signatures being linear combinations of two or more other

614    signatures, for each identified signature, a synthetic sample was created with the

615    pattern of the signature multiplied by 1,000 copy number segments. Further, the

616    synthetic sample was resampled with probabilities $p_{i,f} = d_{i,f} / \sum_{j=1}^{48} d_{j,f}$, where $d_{i,f}$ is

617    the strength of the $i^{\text{th}}$ copy number category in the $f^{\text{th}}$ identified signature. Each

618    resampling was then scanned for activity of all other signatures from the reference

619    set. If a resampled sample can be reconstituted with a cosine similarity >0.95 by 3 or

620    fewer other signatures, the signature used to create the synthetic sample was

621    deemed to be a linear combination of those signatures, and the signature was

622    removed from the global reference set of signatures.

623

624    **Reference set of copy number signatures**

625    Initially 28 pan-cancer copy number signatures were derived from the different

626    SigProfilerExtractor analyses of the 9,873 copy number profiles from SNP

627    microarrays. *In silico* evaluation and manual curation showed that 10 copy number

628    signatures were linear combinations of two or more other signatures. Additionally, 3

629    signatures were deemed to be artefactual due to over-segmentation of copy number

630    profiles. These artefactual signatures were removed from further analyses, as were

631    the samples with any attribution of any of these artefactual signatures (116 samples;

632    1.2% of all TCGA samples). Moreover, samples with >25Mb of homozygous

633    deletions across the genome were removed from downstream analysis (58

634    samples), leaving 9,699 samples for full analysis. Upon signature assignment (see

635     below) 3 of the signatures that were removed due to linear combination were re-

636     extracted within tumour-type specific assignment (cosine similarity=1), suggesting

637     some copy number profiles could not be explained well without these 3 signatures.

638     As a result, these 3 signatures were reintroduced into the compendium of signatures,

639     leaving a total of 19 non-artefactual pan-cancer signatures of copy number

640     alteration.

641

642     CN1-3 form a group of ploidy-associated signatures. CN1 and CN2 display TCN

643     between 2 and 3-4 respectively, with predominantly >40Mb heterozygous segments.

644     CN3 consists of predominantly heterozygous segments of TCN 5-8 with sizes >1Mb.

645

646     CN4-8 form a group of amplicon-associated signatures, that all have segment sizes

647     predominantly between 100kb and 10Mb but with differing TCN or LOH states. CN4

648     consists of a mixture of LOH segments with TCN 1 and heterozygous segments with

649     TCN 3-4. CN5 consists almost entirely of LOH segments with TCN 2. CN6 consists

650     of a mixture of LOH segments with TCN 2 and heterozygous segments with TCN 3-

651     4. CN7 consists of a mixture of heterozygous segments with TCN of 3-4, 5-8 and 9+.

652     CN8 consists of predominantly heterozygous segments with TCN 9+.

653

654     CN9-12 form a group of signatures with considerable LOH components. CN9

655     consists of a mixture of LOH segments with TCN 2 and heterozygous segments with

656     TCN 2, each ranging from 100kb-40Mb. CN10 consists of a mixture of LOH

657     segments with TCN 2 and 3-4 as well as heterozygous segments with TCN 3-4

658     between 100kb and 40Mb. CN11 consists of a mixture of LOH segments with TCN

659     3-4 and heterozygous segments with TCN 5-8, each at predominantly 1-10Mb. CN12

660    consists of mostly LOH segments of TCN 2 with sizes above 100kb and additional

661    heterozygous segments of TCN 3-4 with sizes between 10 and 40Mb.

662

663    CN13-14 form a group of signatures with whole-arm or whole-chromosome scale

664    LOH events. CN13 consists of LOH segments with TCN 2 and heterozygous

665    segments with TCN 3-4, each at >40Mb, while CN14 is similar but with TCN 3-4 and

666    5-8 for LOH and heterozygous segments respectively.

667

668    CN15 has been associated with the tandem duplicator phenotype (**Fig. 4**). This

669    signature consists of LOH segments of TCN 2 and 3-4 as well as heterozygous

670    segments of TCN 3-4 and 5-8, each with segment sizes 1-40Mb.

671

672    CN16-19 originate from unknwon processes and are diverse in their copy number

673    patterns. CN16 consists of predominantly heterozygous segments of TCN 4-8 at

674    >1Mb, but with appreciable contributions of LOH segments with TCN 3-4 at >1Mb

675    and heterozygous segments with TCN 9+ at >100kb. CN17 consists of segments

676    between 100kb and 40Mb that are heterozygous with TCN 3-4 or less commonly

677    LOH with TCN 1 or 2. CN18 consists of predominantly heterozygous segments with

678    TCN 3-4 at 100kb-40Mb with some heterozygous segments of TCN 3-4 at 100kb-

679    10Mb. CN19 consists of heterozygous segments with TCN 2 at >1Mb and many

680    heterozygous segments with TCN 3-4 at 100kb-1Mb.

681

682    **Assignment of copy number signatures to individual cancer samples**

683    The global reference set of copy number signatures was used to assign an activity

684    for each signature to each of 9,873 examined samples using the decomposition

685   module of the SigProfilerExtractor[26]. For the assignment, the information of the *de*

686   *novo* signature and their activities assigned to each sample were used to implement

687   the decomposition module with default parameters except for the NNLS addition

688   penalty (*nnls_add_penalty*) which was set to 0.1, the NNLS removal penalty

689   (*nnls_remove_penalty*) which was set to 0.01, and the initial removal penalty

690   (*initial_remove_penalty*) which was set to 0.05. Signatures were assigned to

691   samples in both tumour-specific evaluations and in a pan-cancer evaluation. As

692   previously done[60], the signature attributions from either tumour-specific or pan-

693   cancer evaluations that gave the best cosine similarity between the input sample

694   vector and the reconstructed sample vector were used as the attributions for that

695   sample in all subsequent analyses.

696

697   **Copy number signatured derived from whole-genome and exome sequencing**

698   **data**

699   A set of samples from TCGA with both SNP-array and exome sequencing data were

700   selected (*n*=282). Copy number profiles were generated from the exome sequencing

701   data using ASCAT across all of the dbSNP common SNP positions with a

702   segmentation penalty ranging from 20 to 140. Signatures were re-extracted for these

703   282 samples from both the SNP-array derived copy number profiles and the exome-

704   derived copy number profiles, and the resulting signatures were compared.

705

706   For whole-genome sequencing data, we examined 512 whole-genome sequenced

707   samples from the PCAWG project overlapping with TCGA samples with microarray

708   data. Copy number profiles from whole-genome sequencing data were generated

709   using ASCAT across the SNP6 positions, with a segmentation penalty ranging from

710 20 to 120. Signatures were extracted for samples with both SNP6 microarray derived

711 copy number profiles and the WGS derived copy number profiles, and the extracted

712 signatures were compared. In all cases, segmentation penalty of 70 gave the best

713 concordance for both copy number profiles and extracted copy number signatures

714 based on SNP6 microarray, whole-genome sequencing, and whole-exome

715 sequencing data.

716

717 **Mapping copy number signatures to the landscapes of cancer genomes**

718 Given the original copy number profiles, the identified signature matrix of $c$ copy

719 number classes by $f$ signatures, and the signature activity matrix of $s$ samples by $f$

720 signatures, it is then possible to map signatures to the genomic landscape for each

721 cancer sample. The probability of each copy number class, $\boldsymbol{c}$, having originated from

722 each signature, $\boldsymbol{i}$ from a total of $\boldsymbol{I}$ signatures, in a sample $\boldsymbol{j}$ can be defined as:

723
$$m_{i,j,c} = \frac{f_{c,i}e_{i,j}l_j}{\sum_{k=1}^{I} f_{c,k}e_{k,j}l_j},$$

724 where $\boldsymbol{f}$ is the normalised signature matrix, $\boldsymbol{e}$ is the normalized attribution matrix,

725 and $\boldsymbol{l}$ is a matrix of the number of segments in the copy number profile of each

726 sample. The likelihood of each signature contributing to a given genomic window,

727 here taken as each chromosome, is then the sum of copy number class probabilities

728 for each segment in that window:

729
$$p_{i,j,w} = \sum_{x=1}^{l_{j,w}} m_{i,j,c_x}$$

730 Once these chromosome likelihoods have been calculated, the individual segments

731 in a chromosome are assigned to their maximum likelihood signature. Once copy

732 number signatures have been mapped to the genome at a segment level, it is

733 possible to interrogate the recurrence of signatures across the genome for a given

33

734    set of copy number profiles. To do this, the genome is binned into 1Mb tiled

735    windows. Within each window, the number of samples with a segment of a given

736    copy number signature that overlaps the window is computed. This is repeated for

737    each signature in each window.

738

**739    Associations between copy number signatures and events defined by genomic**

**740    region**

741    Localised events (chromothripsis[39] and amplicon structure[38]) identified using WGS

742    data were associated with mapped copy number signatures from TCGA for all

743    available matching samples (chromothripsis $n$=657; amplicon $n$=1703). Each

744    segment in every sample was categorised as overlapping or non-overlapping of a

745    localized event. For each copy number signature, the association was then tested

746    using a two-sided Fisher's exact test on a contingency table of segments categorized

747    as overlapping or non-overlapping of a localized event and assigned to or not

748    assigned to the given copy number signature, across all samples. Multiple-testing

749    correction was performed using the Benjamini-Hochberg method.

750

**751    Genome doubled copy number signatures**

752    With the copy number categories being defined as 0, 1, 2, 3-4, 5-8, and 9+, it is

753    possible to artificially 'genome double' any copy number category, other than 0, by

754    assigning it to the next highest copy number category. In this way we artificially

755    'genome doubled' each signature by assigning the count for each copy number class

756    to its next highest copy number class. First, the copy number 1 class is assigned a

757    count of 0, then each copy number class is assigned the count of the preceding copy

758    number class. For example, copy number class of 2 is assigned to the previous copy

759    number class of 1, 3-4 assigned previous 2, *etc.*, until finally the copy number 9+

760    class is assigned a count that is the sum of the previous copy number 5-8 class and

761    9+ class. During this conversion, LOH and size categories are retained, so that the

762    only shift is in copy number. Having performed this conversion, cosine similarities

763    between the artificially 'genome doubled' signatures and the original signatures were

764    calculated. Any genome-doubled and original signature pair that had a cosine

765    similarity >0.85 was considered to contain a pair of signatures with analogous copy

766    number patterns distinguished only by their genome doubling status.

767

768    **Associations between copy number signatures and ploidy**

769    Ploidy for each copy number profile was calculated as the relative length weighted

770    sum of total copy number across a sample. The proportions of the genome that

771    displayed LOH (pLOH) were also calculated. Samples with a ploidy above -

772    3/2*pLOH+3, meaning an LOH-adjusted ploidy of 3 or greater were deemed to be

773    genome doubled samples, while samples with a ploidy above -5/2*pLOH+5,

774    meaning an LOH-adjusted ploidy of 5 or greater, were deemed to be twice genome

775    doubled samples. All other samples were considered as non-genome doubled

776    samples. Each signature (CN1-19) was associated with each genome doubling

777    category (GDx0, GDx1, and GDx2) using a one-sided Fisher's exact test on a

778    contingency table with samples categorized by whether the samples have >0.05

779    attribution to the given copy number signature or not, and whether the sample has

780    the given genome doubled category or not. All p-values were corrected for multiple

781    hypothesis testing using the Benjamini-Hochberg method.

782

**Associations between copy number signatures and known cancer risk factors**

Associations between attributions of copy number signatures and attributions of single-base substitutions, indels, and doublet base signature exposures[14] were performed using Kendall's rank correlation. Only the significant associations found in both cancer-type specific and pan-cancer analysis were reported. For the cancer risk association analyses, copy number signatures were associated with gender[61], tobacco smoking[18], and alcohol drinking status[62]. For each copy number signature, the association was conducted using a two-sided Fisher's exact test on a contingency table of a clinical feature categorized as present or absent and assigned to or not assigned to the given copy number signature across all samples. All p-values were corrected for multiple hypothesis testing using the Benjamini-Hochberg method.

Associations between copy number signature attribution (binarized to present or absent) and the tandem duplicator phenotype (also binarized to present or absent)[35] were performed using a two-sided Fisher's exact test ($n$=882). This was performed for each copy number signature separately. All p-values were corrected for multiple hypothesis testing using the Benjamini-Hochberg method and only associations with q<0.05 were reported.

Associations between copy number signature attribution (binarized to present or absent) and driver gene SNV/indel mutation status[63] were performed within tumour types using a two-sided Fisher's exact test ($n$=6,543 across all cancer types). This was performed for all copy number signature/gene combinations for which the gene was mutated in the given cancer type and the copy number signature was observed

36

808    in the given cancer type. All p-values were corrected for multiple hypothesis testing

809    using the Benjamini-Hochberg method and only associations with both q<0.05 and

810    |log2(OR)|>1 were reported.

811

812    Driver copy number alterations of COSMIC cancer gene census genes[55] were

813    defined as: *(i)* homozygous deletion (CN={0,0}) of genes listed as deleted (D) in

814    COSMIC mutation types; or *(ii)* amplification (CN>2*ploidy+1) of genes listed as

815    amplified (A) in COSMIC mutation types. Associations were then performed on copy

816    number driver alterations for SNV/indel driver gene alterations as above (*n*=9,699

817    across all cancer types).

818

819    The diversity of copy number signatures, as defined by Shannon's diversity index,

820    was associated with both SNV/indel and copy number driver gene mutations using a

821    logistic regression model with binary diversity {>0, =0} as the dependent variable,

822    and tumour type and gene mutation status as independent variables. LGG was taken

823    as the reference tumour type. Only driver genes with >250 mutant samples in the

824    dataset were included in the model.

825

826    Associations between copy number signature attribution (binarized to present or

827    absent) and age at diagnosis (binarized to above or below median separately for

828    each cancer type) were performed within cancer types using a two-sided Fisher's

829    exact test (*n*=8,841 across all cancer types). All p-values were corrected for multiple

830    hypothesis testing using the Benjamini-Hochberg method and only associations with

831    both q<0.05 and |log2(OR)|>1 were reported.

832

833

**Copy number signatures and defective homologous recombination**

835 Signatures were tested for enrichment in tumour types using one-sided Mann-

836 Whitney tests of signature attribution in a given tumour type versus all other tumour

837 types. This was performed for all signature and tumour combinations. All p-values

838 were corrected for multiple hypothesis testing using the Benjamini-Hochberg

839 method.

840

841 Core homologous recombination (HR) repair pathway member genes were chosen

842 to interrogate: *BRCA1*, *BRCA2*, *RAD51C*, *PALB2*[64,65]. Copy number alterations

843 across these genes were identified based on ASCAT copy number profiles for

844 homozygous deletions (*i.e.,* CN={0, 0}) and LOH (*i.e.,* CN={>0, 0}). Somatic SNVs

845 and indels were taken from Ref. [63]. Pathogenic germline variants in *BRCA1* and

846 *BRCA2* were taken from Ref. [66]. Samples were deemed as bi-allelically mutated for

847 the HR pathway if homozygously deleted (HD) or if >1 of any of the other classes of

848 alteration were present within any of the HR pathway genes. Mono-allelic loss was

849 defined as 1 of any of the non-HD alterations within any of the HR pathway genes.

850 Wildtype was defined as no alterations in any HR pathway genes. The associations

851 between HR pathway status and CN15 were then restricted to only breast (*n*=589),

852 ovarian (*n*=309), and pan-cancer (*n*=4,919). Two-sided fisher's exact tests were

853 performed between wild-type and mono-allelic samples, between wild-type and bi-

854 allelic samples, and between mono-allelic and bi-allelic HR pathway status samples.

855 All p-values were corrected for multiple hypothesis testing using the Benjamini-

856 Hochberg method.

857

**Copy number signatures associated with changes of overall survival**

Survival data for 11,160 TCGA patients were obtained from the TCGA Clinical data Resource R package[67]. Univariate disease specific survival analysis for signatures was performed using a log-rank test and Kaplan-Meier curves in R, with groups being unattributed (attribution=0) and attributed (attribution>0) for each signature separately, or for summed attributions of a set of signatures (*e.g.*, amplicon signatures).

Multivariate disease-specific survival analysis was performed using the Cox's proportional hazards model in R with Boolean attributed/non-attributed variables for each copy number signature and tumour type as covariates. To account for potential violations of Cox's model's proportional hazards assumption, we also conducted the same analysis using the accelerated failure time model with the Weibull distribution using the flexsurvreg function in R. All p-values were corrected for multiple hypothesis testing using the Benjamini-Hochberg method.

**Simulating copy number profiles**

*Simulation framework:* Genomes were initialized as 23 pairs of individual chromosomes, with lengths corresponding to those seen in the human genome, where the 23$^{rd}$ pair could be either *X, X* or *X, Y*. Each chromosome was initialized as a data table with chromosome (1-22, X, Y), start position, end position, and allele (either A or B). Genomic events were recorded as altering one of these data tables in the appropriate way, adding or removing segments as necessary. Gains and losses: The $\log_{10}$(size) of sub-chromosomal gains were drawn from a Gaussian mixture with components:

883 $$\mathbf{N}(\mu=5.961351, \sigma^2=0.4199448),$$

884 $$\mathbf{N}(\mu=7.786183, \sigma^2=0.1068539),$$

885 at proportions $p_1=0.7360366$ and $p_2=1-p_1$. The $\log_{10}$(size) of sub-chromosomal

886 losses were drawn from a gaussian mixture with components:

887 $$\mathbf{N}(\mu=6.188331, \sigma^2= 0.5686788),$$

888 $$\mathbf{N}(\mu=7.588125, \sigma^2= 0.1326166),$$

889 at proportions $p_1=0.6472512$ and $p_2=1-p_1$. The parameters for the various

890 distributions were estimated from samples in TCGA that were predominantly diploid

891 (CN1+CN9 attribution>0.8) from segments that were copy number 1 for the loss

892 distributions, and copy number 3 for the gain distributions. Parameters were

893 estimated using a Gaussian mixture model on the $\log_{10}$(sizes) of the appropriate

894 segments with two components due to the bimodal nature of the segment length

895 distributions.

896

897 First the chromosome on which the gain/loss will occur is randomly sampled with

898 probabilities $1/n$, where $n$ is the number of separate chromosomes in the current

899 genome. The event size, $\lambda$ ,is then drawn from the previously stated multinormal

900 distributions; if an event size greater than the chromosomal size is drawn, then a

901 new size is drawn. The start of the event, $b_1$, is then drawn from a uniform

902 distribution,

903 $$b_1 \sim \mathbf{U}(1, e-\lambda),$$

904 where $e$ is the cumulative length of the chosen chromosome, and the end of the

905 event, $b_2=b_1+\lambda$.

906

40

907 Gains are treated as tandem duplications, so that the gained region is inserted

908 immediately after the start breakpoint. On unaltered chromosome, this will alter the

909 chromosome from a single segment with start=1 and end=$e$ to a chromosome with

910 four segments, with starts=[1,$b_1$+1,$b_1$+1,$b_2$+1] and ends=[$b_1$,$b_2$,$b_2$,$e$], each with the

911 chosen chromosome identity and allele; note that this will eventually lead to a copy

912 number profile with 3 segments with starts==[1,$b_1$+1,$b_2$+1] and ends=[$b_1$,$b_2$,$e$]. A loss

913 will instead lead to a chromosome with two segments with starts=[1,$b_2$] and

914 ends=[$b_1$,$e$].

915

916 *Simulating chromothripsis:* For chromothriptic events, the $\log_{10}$(number of segments)

917 for the resulting chromosome is drawn from a normal distribution:

918         $n \sim \mathbf{N}(\mu=1.3, \sigma=0.3),$

919 while the $\log_{10}$(length) of segments are drawn from a normal distribution

920         $\lambda \sim \mathbf{N}(\mu=6, \sigma=0.7),$

921 and the start of the chromothriptic event is drawn from a uniform distribution:

922        $\mathbf{U}(1, e - \sum_{1}^{n} \lambda_n),$

923 where $e$ is the size of the chromosome. The parameters for the distributions were

924 chosen to match the empirical distributions observed in TCGA chromosomes that

925 were called as chromothriptic in the PCAWG dataset.

926

927 The breakpoints of the chromothriptic event, [$b_1$,…,$b_{n-1}$], are then the cumulative

928 sums of the segment sizes, apart from the first breakpoint which is 1. The

929 chromosome is then broken into $n$ segments by their cumulative lengths, defined by

930 the breakpoints. Whether to lose a segment is drawn from a binomial distribution:

931        $\delta_x \sim \mathbf{Binom}(1,0.5).$

932    All segments were removed where $\delta_x$=1. The remaining segments were then

933    randomly reversed if:

934                                    $\rho_x$~**Binom**(1,0.5)=1.

935    Lastly, the remaining segments were resampled without replacement so that their

936    order is randomized, and are then concatenated together. The chromothriptic

937    chromosome replaces the original chromosome that it originates from.

938

939    *Genome doubling and chromosomal gains/losses:* All chromosomes in the set of

940    chromosomes are duplicated to simulate genome doubling. For chromosomal gains,

941    a single chromosome is duplicated, whereas for chromosomal losses a single

942    chromosome is removed.

943

944    *Calculating copy number:* Once an assortment of chromosomes has been simulated

945    from a mixture of the previously described processes, the combined copy number

946    across all derivative chromosomes must be calculated across the reference genome.

947    For each reference chromosome, *x*, all segments across the derivative

948    chromosomes that derive from *x* are collated, and the breakpoints across *x* are

949    defined as the ordered unique set of start or end positions of those segments. Then

950    the copy number for segment $i_x$, is calculated for each allele separately; the A allele

951    copy number is the count of A allele segments in all derivative chromosomes that

952    overlap the segment defined between $b_{i,x}$ and $b_{i+1,x}$, and similar for the B allele copy

953    number. Combined across all reference chromosomes, this gives an allele-specific

954    copy number profile.

955

42

956    *Combinations of simulations:* The following simulations were performed, for 100

957    samples each:

958        • CINx10 – 10 random gain or loss events.

959        • CINx50 – 50 random gain or loss events.

960        • CINx10->WGD – 10 random gain or loss events, followed by WGD.

961        • CINx50->WGD – 50 random gain or loss events, followed by WGD.

962        • CINx5->WGD->CINx50 - 5 random gain or loss events, followed by WGD,

963          followed by 50 random gain or loss events.

964        • CINx5->WGD->CINx25->WGD->CINx25 - 5 random gain or loss events,

965          followed by WGD, followed by 25 random gain or loss events, followed by

966          WGD, followed by 25 random gain or loss events.

967        • Chromo. – Chromothripsis of a random chromosome.

968        • Chromo.->WGD – Chromothripsis of a random chromosome, followed by

969          WGD.

970        • Chromo.->Amp. – Chromothripsis of a random chromosome, followed by

971          chromosomal gain of the derivative chromothriptic chromosome.

972        • Chromo.->Amp.->WGD - Chromothripsis of a random chromosome, followed

973          by chromosomal gain of the derivative chromothriptic chromosome, followed

974          by WGD.

975        • Chromo.->Amp.x5->WGD. Chromothripsis of a random chromosome,

976          followed by chromosomal gain of the derivative chromothriptic chromosome

977          five times, followed by WGD.

978    For random gain/loss events, a binomial draw was used to decide whether a gain or

979    loss occurred, with $p_{gain}$=0.4.

980

981

**REFERENCE**

1.  Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: The next generation. *Cell* **144**, 646–674 (2011).
2.  Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates with markers of immune evasion and with reduced response to immunotherapy. *Science (80-. ).* **355**, (2017).
3.  Taylor, A. M. *et al.* Genomic and Functional Approaches to Understanding Cancer Aneuploidy. *Cancer Cell* **33**, 676-689 e3 (2018).
4.  Ben-David, U. & Amon, A. Context is everything: aneuploidy in cancer. *Nature Reviews Genetics* **21**, 44–62 (2020).
5.  Rajagopalan, H. & Lengauer, C. Aneuploidy and cancer. *Nature* **432**, 338–341 (2004).
6.  Beroukhim, R. *et al.* The landscape of somatic copy-number alteration across human cancers. *Nature* **463**, 899–905 (2010).
7.  Sansregret, L. & Swanton, C. The role of aneuploidy in cancer evolution. *Cold Spring Harbor Perspectives in Medicine* **7**, (2017).
8.  Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of advanced cancers. *Nat. Genet.* **50**, 1189–1195 (2018).
9.  Stephens, P. J. *et al.* Massive genomic rearrangement acquired in a single catastrophic event during cancer development. *Cell* **144**, 27–40 (2011).
10. Bolhaqueiro, A. C. F. *et al.* Ongoing chromosomal instability and karyotype evolution in human colorectal cancer organoids. *Nat. Genet.* **51**, 824–834 (2019).
11. Shoshani, O. *et al.* Chromothripsis drives the evolution of gene amplification in cancer. *Nature* 1–5 (2020). doi:10.1038/s41586-020-03064-z
12. Umbreit, N. T. *et al.* Mechanisms generating cancer genome complexity from a single cell division error. *Science (80-. ).* **368**, (2020).
13. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer (vol 500, pg 415, 2013). *Nature* **502**, (2013).
14. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature* **578**, 94–101 (2020).
15. Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat Genet* **47**, 1402–1407 (2015).
16. Gulhan, D. C., Lee, J. J., Melloni, G. E. M., Cortes-Ciriano, I. & Park, P. J. Detecting the mutational signature of homologous recombination deficiency in clinical samples. *Nat Genet* **51**, 912–919 (2019).
17. Behjati, S. *et al.* Mutational signatures of ionizing radiation in second malignancies. *Nat. Commun.* **7**, (2016).
18. Alexandrov, L. B. *et al.* Mutational signatures associated with tobacco smoking in human cancer. *Science* **354**, 618–622 (2016).
19. Haradhvala, N. J. *et al.* Distinct mutational signatures characterize concurrent loss of polymerase proofreading and mismatch repair. *Nat Commun* **9**, 1746 (2018).
20. Kucab, J. E. *et al.* A Compendium of Mutational Signatures of Environmental Agents. *Cell* **177**, 821-836 e16 (2019).
21. Meier, B. *et al.* Mutational signatures of DNA mismatch repair deficiency in C. elegans and human cancers. *Genome Res* **28**, 666–675 (2018).
22. Pich, O. *et al.* The mutational footprints of cancer therapies. *Nat Genet* **51**, 1732–1740 (2019).
23. Steele, C. D. *et al.* Undifferentiated Sarcomas Develop through Distinct

1032            Evolutionary Pathways. *Cancer Cell* **35**, 441-456.e8 (2019).

1033   24.   Macintyre, G. *et al.* Copy number signatures and mutational processes in
1034            ovarian carcinoma. *Nat Genet* **50**, 1262–1270 (2018).

1035   25.   Pladsen, A. V. *et al.* DNA copy number motifs are strong and independent
1036            predictors of survival in breast cancer. *Commun. Biol.* **3**, 1–9 (2020).

1037   26.   Ashiqul Islam, S. M. *et al.* Uncovering novel mutational signatures by de novo
1038            extraction with. *bioRxiv* 2020.12.13.422570 (2020).
1039            doi:10.1101/2020.12.13.422570

1040   27.   Alexandrov, L. B. & Stratton, M. R. Mutational signatures: The patterns of
1041            somatic mutations hidden in cancer genomes. *Current Opinion in Genetics and*
1042            *Development* **24**, 52–60 (2014).

1043   28.   Koh, G., Zou, X. & Nik-Zainal, S. Mutational signatures: Experimental design
1044            and analytical framework. *Genome Biology* **21**, 37 (2020).

1045   29.   López, S. *et al.* Interplay between whole-genome doubling and the
1046            accumulation of deleterious alterations in cancer evolution. *Nat. Genet.* **52**,
1047            283–293 (2020).

1048   30.   Davoli, T. *et al.* Cumulative haploinsufficiency and triplosensitivity drive
1049            aneuploidy patterns and shape the cancer genome. *Cell* **155**, 948–962 (2013).

1050   31.   Mertens, F., Johansson, B., Höglund, M. & Mitelman, F. Chromosomal
1051            Imbalance Maps of Malignant Solid Tumors: A Cytogenetic Survey of 3185
1052            Neoplasms. *Cancer Res.* **57**, 2765–2780 (1997).

1053   32.   Kim, H. *et al.* Extrachromosomal DNA is associated with oncogene
1054            amplification and poor outcome across multiple cancers. *Nat. Genet.* **52**, 891–
1055            897 (2020).

1056   33.   Zheng, S. *et al.* Comprehensive Pan-Genomic Characterization of
1057            Adrenocortical Carcinoma. *Cancer Cell* **30**, 363 (2016).

1058   34.   Davis, C. F. *et al.* The somatic genomic landscape of chromophobe renal cell
1059            carcinoma. *Cancer Cell* **26**, 319–330 (2014).

1060   35.   Menghi, F. *et al.* The Tandem Duplicator Phenotype Is a Prevalent Genome-
1061            Wide Cancer Configuration Driven by Distinct Gene Mutations. *Cancer Cell* **34**,
1062            197-210.e5 (2018).

1063   36.   Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of
1064            advanced cancers. *Nat Genet* **50**, 1189–1195 (2018).

1065   37.   Lo, A. W. I. *et al.* DNA amplification by breakage/fusion/bridge cycles initiated
1066            by spontaneous telomere loss in a human cancer cell line. *Neoplasia* **4**, 531–
1067            538 (2002).

1068   38.   Deshpande, V. *et al.* Exploring the landscape of focal amplifications in cancer
1069            using AmpliconArchitect. *Nat. Commun.* **10**, 1–14 (2019).

1070   39.   Consortium, I. T. P.-C. A. of W. G. Pan-cancer analysis of whole genomes.
1071            *Nature* **578**, 82–93 (2020).

1072   40.   Behjati, S. *et al.* Recurrent mutation of IGF signalling genes and distinct
1073            patterns of genomic rearrangement in osteosarcoma. *Nat. Commun.* **8**, (2017).

1074   41.   Cortes-Ciriano, I. *et al.* Comprehensive analysis of chromothripsis in 2,658
1075            human cancers using whole-genome sequencing. *Nat Genet* (2020).
1076            doi:10.1038/s41588-019-0576-7

1077   42.   Furgason, J. M. *et al.* Whole genome sequence analysis links chromothripsis
1078            to EGFR, MDM2, MDM4, and CDK4 amplification in glioblastoma.
1079            *Oncoscience* **2**, 618–628 (2015).

1080   43.   Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular
1081            Classification of 10,000 Tumors from 33 Types of Cancer. *Cell* **173**, 291-304

1082    e6 (2018).

1083  44.  Knudson, A. G. Hereditary cancer: two hits revisited. *J Cancer Res Clin Oncol*
1084      **122**, 135–140 (1996).

1085  45.  Scarpa, A. *et al.* Whole-genome landscape of pancreatic neuroendocrine
1086      tumours. *Nature* **543**, 65–71 (2017).

1087  46.  Ricketts, C. J. *et al.* The Cancer Genome Atlas Comprehensive Molecular
1088      Characterization of Renal Cell Carcinoma. *Cell Rep.* **23**, 313-326.e5 (2018).

1089  47.  Menghi, F. *et al.* The tandem duplicator phenotype as a distinct genomic
1090      configuration in cancer. *Proc. Natl. Acad. Sci. U. S. A.* **113**, E2373–E2382
1091      (2016).

1092  48.  McBride, D. J. *et al.* Tandem duplication of chromosomal segments is common
1093      in ovarian and breast cancer genomes. *J. Pathol.* **227**, 446–455 (2012).

1094  49.  TCGA. Comprehensive and Integrated Genomic Characterization of Adult Soft
1095      Tissue Sarcomas. *Cell* **171**, 950-965 e28 (2017).

1096  50.  Sakofsky, C. J. *et al.* Repair of multiple simultaneous double-strand breaks
1097      causes bursts of genome-wide clustered hypermutation. *PLoS Biol.* **17**,
1098      e3000464 (2019).

1099  51.  Pfister, K. *et al.* Identification of Drivers of Aneuploidy in Breast Tumors. *Cell*
1100      *Rep.* (2018). doi:10.1016/j.celrep.2018.04.102

1101  52.  Schjølberg, A. R., Clausen, O. P. F., Burum-Auensen, E. & De Angelis, P. M.
1102      Aneuploidy is associated with TP53 expression but not with BRCA1 or TERT
1103      expression in sporadic colorectal cancer. *Anticancer Res.* (2009).

1104  53.  Cazzola, A. *et al.* TP53 deficiency permits chromosome abnormalities and
1105      karyotype heterogeneity in acute myeloid leukemia. *Leukemia* (2019).
1106      doi:10.1038/s41375-019-0550-5

1107  54.  Thorsson, V. *et al.* The Immune Landscape of Cancer. *Immunity* **48**, 812-
1108      830.e14 (2018).

1109  55.  Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer.
1110      *Nucleic Acids Res.* **47**, D941–D947 (2019).

1111  56.  De Jonge, M. M. *et al.* Frequent homologous recombination deficiency in high-
1112      grade endometrial carcinomas. *Clin. Cancer Res.* **25**, 1087–1097 (2019).

1113  57.  Davies, H. *et al.* HRDetect is a predictor of BRCA1 and BRCA2 deficiency
1114      based on mutational signatures. *Nat. Med.* **23**, 517-+ (2017).

1115  58.  Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proc Natl*
1116      *Acad Sci U S A* **107**, 16910–16915 (2010).

1117  59.  Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer.
1118      *Nature* **500**, 415–421 (2013).

1119  60.  Health, N. *et al.* Signatures of mutational processes in human cancer. *Nature*
1120      1–108 (2013). doi:10.1038/nature

1121  61.  Liu, J. *et al.* An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive
1122      High-Quality Survival Outcome Analytics. *Cell* **173**, 400-416.e11 (2018).

1123  62.  Lowy, D. R., Kibbe, W. A., Ph, D., Staudt, L. M. & Ph, D. New engla nd journal.
1124      1109–1112 (2016).

1125  63.  Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic
1126      Tissues. *Cell* **171**, 1029-1041 e21 (2017).

1127  64.  Knijnenburg, T. A. *et al.* Genomic and Molecular Landscape of DNA Damage
1128      Repair Deficiency across The Cancer Genome Atlas. *Cell Rep* **23**, 239-254 e6
1129      (2018).

1130  65.  Nguyen, L., W. M. Martens, J., Van Hoeck, A. & Cuppen, E. Pan-cancer
1131      landscape of homologous recombination deficiency. *Nat. Commun.* **11**, 1–12

1132       (2020).
1133   66.   Yost, S., Ruark, E., Alexandrov, L. B. & Rahman, N. Insights into BRCA
1134       Cancer Predisposition from Integrated Germline and Somatic Analyses in 7632
1135       Cancers. *JNCI Cancer Spectr.* **3**, (2019).
1136   67.   Liu, J. *et al.* An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive
1137       High-Quality Survival Outcome Analytics. *Cell* **173**, 400-416.e11 (2018).
1138
1139
1140

1141
1142