

# 1 Pan-cancer whole genome analyses of metastatic solid tumors

2 Peter Priestley<sup>1,2,\*#</sup>, Jonathan Baber<sup>1,2,\*</sup>, Martijn P. Lolkema<sup>3,4</sup>, Neeltje Steeghs<sup>3,5</sup>, Ewart de Bruijn<sup>1</sup>,  
3 Korneel Duyvesteyn<sup>1</sup>, Susan Haidari<sup>1,3</sup>, Arne van Hoeck<sup>6</sup>, Wendy Onstenk<sup>1,3,4</sup>, Paul Roepman<sup>1</sup>,  
4 Charles Shale<sup>2</sup>, Mircea Voda<sup>1</sup>, Haiko J. Bloemendaal<sup>7</sup>, Vivianne C.G. Tjan-Heijnen<sup>8</sup>, Carla M.L. van  
5 Herpen<sup>9</sup>, Mariette Labots<sup>10</sup>, Petronella O. Witteveen<sup>11</sup>, Egbert F. Smit<sup>3,5</sup>, Stefan Sleijfer<sup>3,4</sup>, Emile E.  
6 Voest<sup>3,5</sup>, Edwin Cuppen<sup>1,3,6,#</sup>

7

8 <sup>1</sup> Hartwig Medical Foundation, Science Park 408, Amsterdam, The Netherlands

9 <sup>2</sup> Hartwig Medical Foundation Australia, Sydney, Australia

10 <sup>3</sup> Center for Personalized Cancer Treatment, The Netherlands

11 <sup>4</sup> Erasmus MC Cancer Institute, Doctor Molewaterplein 40, Rotterdam, The Netherlands

12 <sup>5</sup> Netherlands Cancer Institute/Antoni van Leeuwenhoekhuis, Plesmanlaan 121, Amsterdam, The Netherlands

13 <sup>6</sup> Center for Molecular Medicine and Oncode Institute, University Medical Center Utrecht, Heidelberglaan 100,  
14 Utrecht, The Netherlands

15 <sup>7</sup> Meander Medisch Centrum, Maatweg 3, Amersfoort, The Netherlands

16 <sup>8</sup> Maastricht University Medical Center, P. Debyelaan 25, Maastricht, The Netherlands

17 <sup>9</sup> Radboud University Medical Center, Geert Grooteplein Zuid 10, Nijmegen, The Netherlands

18 <sup>10</sup> VU Medical Center, De Boelelaan 1117, Amsterdam, The Netherlands

19 <sup>11</sup> Cancer Center, University Medical Center Utrecht, Heidelberglaan 100, Utrecht, The Netherlands

20

21 \* shared first author

22 # corresponding authors: p.priestley@hartwigmedicalfoundation.nl, e.cuppen@hartwigmedicalfoundation.nl

23

## Abstract

24

Metastatic cancer is one of the major causes of death and associated with poor treatment efficiency. A better understanding of the characteristics of late stage cancer is required to help tailor personalised treatment, reduce overtreatment and improve outcomes. Here we describe the largest pan-cancer study of metastatic solid tumor genomes, including 2,520 whole genome-sequenced tumor-normal pairs, analyzed at a median depth of 106x and 38x, respectively, and surveying over 70 million somatic variants. Metastatic lesions were found to be very diverse, with mutation characteristics reflecting those of the primary tumor types, although with high rates of whole genome duplication events (56%). Metastatic lesions are relatively homogeneous with the vast majority (96%) of driver mutations being clonal and up to 80% of tumor suppressor genes bi-allelically inactivated through different mutational mechanisms. For 62% of all patients, genetic variants that may be associated with outcome of approved or experimental therapies were detected. These actionable events were distributed over the various mutation types (single and multiple nucleotide variants, insertions and deletions, copy number alterations and structural variants) underlining the importance of comprehensive genomic tumor profiling for cancer precision medicine for advanced cancer treatment.

40

## Introduction

41

Metastatic cancer is one of the leading causes of death globally and is a major burden for society despite the availability of an increasing number of (targeted) drugs. Health care costs associated with treatment of metastatic disease are increasing rapidly due to the high cost of novel targeted treatments and immunotherapy, while many patients do not benefit from these approaches with inevitable adverse effects for most patients. Metastatic cancer therefore poses a major challenge for society to balance between individual and societal treatment responsibilities. As cancer genomes

47 evolve over time, both in the highly heterogeneous primary tumor mass and as disseminated  
48 metastatic cells<sup>1,2</sup>, a better understanding of metastatic cancer genomes is crucial to further improve  
49 on tailoring treatment for late stage cancers.

50 In recent years, several large-scale whole genome sequencing (WGS) analysis efforts such  
51 as TCGA and ICGC have yielded valuable insights in the diversity of the molecular processes driving  
52 different types of adult<sup>3,4</sup> and pediatric<sup>5,6</sup> cancer and have fueled the promises of genome-driven  
53 oncology care<sup>7</sup>. However, these analyses were primarily done on primary tumor material whereas  
54 metastatic cancers, which cause the bulk of the disease burden and 90% of all cancer deaths, have  
55 been less comprehensively studied at the whole genome level, with previous efforts focusing on  
56 tumor-specific cohorts<sup>8-10</sup> or at a targeted gene panel<sup>11</sup> or exome level<sup>12</sup>.

57 Here we describe the first large-scale pan-cancer whole-genome landscape of metastatic  
58 cancers based on the Hartwig Medical Foundation (HMF) cohort of 2,520 paired tumor and normal  
59 genomes from 2,405 patients. The samples have been collected prospectively as fresh frozen  
60 biopsies taken from a broad range of metastases ([Extended Data Fig. 1b](#)) and blood controls from  
61 patients with advanced cancer in a clinical study setup coordinated by the Center for Personalized  
62 Cancer Treatment (CPCT) in 41 hospitals in the Netherlands ([Supplementary Table 1](#)). All samples  
63 were paired with standardized clinical information ([Extended Data Fig. 1a](#)). The sample distribution  
64 over age and primary tumor types broadly reflects solid cancer incidence in the Western world,  
65 including rare cancers ([Figure 1a-b](#)).

66 The cohort has been analyzed with uniform and high depth paired-end (2 x 150 bp) whole  
67 genome sequencing with a median depth of 106x for tumor samples and 38x for the blood control  
68 ([Extended Data Fig. 1c](#)). Sequencing data were analyzed for all types of somatic variants using an  
69 optimized bioinformatic pipeline based on open source tools ([Methods](#)). We identified a total of  
70 59,472,629 single nucleotide variants (SNVs), 839,126 multiple nucleotide variants (MNVs),  
71 9,598,205 insertions and deletions (INDELs) and 653,452 structural variants (SVs) ([Supplementary  
72 Table 2](#)). We found that the relative high sequencing depth is important for variant calling sensitivity  
73 as downsampling of the tumor sample coverage to ~53x resulted in an average decrease in sensitivity  
74 of 10% for SNV, 2% for INDEL, 15% for MNV, and 19% for SV ([Extended Data Fig. 2](#)).

75 Here we present a detailed characterization of this unique and comprehensive resource for a  
76 better genomic understanding of metastatic cancer.

## 77 **Mutational landscape of metastatic cancer**

78 We analysed the tumor mutational burden (TMB) of each class of variants per cancer type  
79 based on tissue of origin ([Fig. 1c-h, Supplementary Table 2](#)). In line with previous studies on primary  
80 cancers<sup>13</sup>, we found extensive variation in mutational load of up to 3 orders of magnitude both within  
81 and across cancer types.

82 The median SNV counts per sample were highest in skin, predominantly consisting of  
83 melanoma (44k) and lung (36k) tumors with ten-fold higher SNV counts than sarcomas (4.1k),  
84 neuroendocrine tumors (NET) (3.5k) and mesotheliomas (3.4k). The variation for MNVs was even  
85 greater with lung (median=815) and skin (median=764) tumors having five times the median MNV  
86 counts of any other tumor type. This can be explained by the well-known mutational impact of UV  
87 radiation (CC->TT MNV) and smoking (CC->AA MNV) mutational signatures, respectively ([Fig. 1f](#)).  
88 Although only di-nucleotide substitutions are typically reported as MNVs, 10.7% of the MNVs involve  
89 three nucleotides and 0.6% had four or more nucleotides affected.

90 INDEL counts were typically ten-fold lower than SNVs, with a lower relative rate for skin and  
91 lung cancers ([Fig. 1d, Extended Data Fig. 3](#)). Genome-wide analysis of INDELs at microsatellite loci  
92 identified 60 samples with microsatellite instability (MSI) ([Supplementary Table 2](#)), representing 2.4%  
93 of all tumors. The highest rates of MSI were observed in central nervous system (CNS) (9.4%), uterus  
94 (9.0%) and prostate (6.1%) tumors. For metastatic colorectal cancer lesions we found an MSI  
95 frequency of only 4.0%, which is lower than reported for primary colorectal cancer, and in line with  
96 better prognosis for patients with localized MSI colorectal cancer, which less often metastasizes<sup>14</sup>.

97 Remarkably, 67% of all INDELS in the entire cohort were found in the 60 MSI samples, and 85% of all  
98 INDELS in the cohort were found in microsatellites or short tandem repeats. Only 0.33% of INDELS  
99 (32k, ~1% of non-microsatellite INDELS) were found in coding sequences, of which the majority (88%)  
100 had a predicted high impact by affecting the open reading frame of the gene.

101 The median rate of SVs across the cohort was 193 per tumor, with the highest median counts  
102 observed in ovary (415) and esophageal (379) tumors, and the lowest in kidney tumors (71) and NET  
103 (56) (Fig. 1h, Supplementary Table 2). Simple deletions were the most commonly observed SV  
104 subtype (33% of all SVs) and were the most prevalent in every cancer type except stomach and  
105 esophageal tumors which were highly enriched in translocations.

## 106 Copy number alteration landscape of metastatic cancer

107 Copy number alterations (CNAs) are important hallmarks of tumorigenesis<sup>15</sup>. Pan-cancer, the  
108 most highly amplified regions in our metastatic cancer cohort contain the established oncogenes such  
109 as EGFR, CCNE1, CCND1 and MDM2 (Fig. 2a). Chromosomal arms 1q, 5p, 8q and 20q are also  
110 highly enriched in moderate amplification across the cohort each affecting >20% of all samples. For  
111 the amplifications of 5p and 8q this is likely related to the common amplification targets of TERT and  
112 MYC, respectively. However, the targets of the amplifications on 1q, predominantly found in breast  
113 cancers (>50% of samples) and amplifications on 20q, predominantly found in colorectal cancers  
114 (>65% of samples), are less clear.

115 We identified some intriguing patterns of recurrent loss of heterozygosity (LOH) caused by  
116 CNAs. Overall an average of 23% of the autosomal DNA per tumor has LOH. Unsurprisingly, TP53  
117 has the highest LOH recurrence at 67% of samples. Many of the other LOH peaks are also explained  
118 by well-known tumor suppressor genes (TSG). However, several clear LOH peaks are observed  
119 which cannot easily be explained by known TSG selection, such as one on 8p (57% of samples),  
120 which could potentially be the result of the mechanism by which the amplification of 8q, the most  
121 commonly amplified part of the genome, is established.

122 There are remarkable differences in LOH between cancer types (Figure 2, Extended Data  
123 Fig. 4). For instance, we observed LOH events on the 3p arm in 90% of kidney samples<sup>16</sup> and LOH of  
124 the complete chromosome 10 in 72% of CNS tumors (predominantly glioblastoma multiforme<sup>17</sup>). Even  
125 in the case of the TP53 region on chromosome 17, different tumor types display clearly different  
126 patterns of LOH. Ovarian cancers exhibit LOH of the full chromosome 17 in 75% of samples whereas  
127 in prostate cancer, which also has 70% LOH for TP53, this is nearly always caused by highly focal  
128 deletions.

129 Unlike LOH events, homozygous deletions are nearly always restricted to small chromosomal  
130 regions. Not a single example was found in which a complete autosomal arm was homozygously  
131 deleted. Homozygous deletions of genes are surprisingly rare as well: we found only 4,915 autosomal  
132 events (average of two per tumor) where one or multiple consecutive genes are fully or partially  
133 homozygously deleted. In 46% of these events a putative TSG was deleted. The scarcity of  
134 passenger homozygous deletions underlines the fact that despite widespread copy number  
135 alterations in metastatic tumors, the vast majority of genes or gross chromosomal organization likely  
136 remain essential for tumor cell survival. Chromosome Y is a special case and is deleted in 35% of all  
137 male tumor genomes, although this varies significantly between cancer types from 5% in CNS tumors  
138 to as high as 68% in biliary tumors.

139 An extreme form of copy number changes can be caused by whole genome duplication  
140 (WGD). We found WGD events in 56% of all samples ranging between 17% in CNS to 80% in  
141 esophageal tumors (Fig. 2d,e). This is much higher than reported previously for primary tumors (25%-  
142 37%)<sup>18,19</sup> and also higher than estimated from panel-based sequencing analyses of advanced tumors  
143 (30%)<sup>20</sup>. Ploidy levels, in combination with tumor purity data, are essential for correct interpretation of  
144 the measured raw SNV and INDEL frequencies, e.g. to discriminate causal bi-allelic inactivation of  
145 TSG from heterozygous passenger events. Hence determining the WGD status of a tumor is highly  
146 relevant for diagnostic applications. Furthermore, WGD has previously been found to correlate with a

147 greater incidence of cancer recurrence for ovarian cancer<sup>19</sup> and has been associated with poor  
148 prognosis across cancer types, independently of established clinical prognostic factors<sup>20</sup>.

149 **Discovery of novel significantly mutated genes**

150 To identify significantly mutated genes (SMGs) potentially specific for metastatic cancer, we  
151 used the dNdScv approach<sup>21</sup> with strict cutoffs ( $q < 0.01$ ) for the pan-cancer and tumor-type specific  
152 datasets. In addition to reproducing previous results on cancer drivers, a few novel genes were  
153 identified (Extended Data Fig. 5, Supplementary Table 3). In the pan-cancer analyses we found only a  
154 single novel SMG, which was not either present in the curated COSMIC Cancer Gene Census or  
155 found by Martincorena et al<sup>21</sup>. This gene, MLK4 ( $q$ -value = 2e-4), is a mixed lineage kinase that  
156 regulates the JNK/P38 and ERK signaling pathways and has been reported to inhibit tumorigenesis in  
157 colorectal cancer<sup>22</sup>. In addition, in our tumor type-specific analyses, we identified a novel breast  
158 cancer-specific SMG - ZFPM1 (also known as Friend of GATA1 (FOG1),  $q$ -value = 8e-5), a zinc-finger  
159 transcription factor protein without clear links with cancer. Nonetheless, we found six unique  
160 frameshift variants (all inactivated biallelically) and three nonsense variants, which suggests a driver  
161 role for this gene in breast cancer.

162 Our cohort also lends support to some prior SMG findings. In particular, eight significantly  
163 mutated putative TSG in the HMF cohort were also found by Martincorena et al<sup>21</sup> - GPS2 (pan-cancer,  
164  $q = 1e-5$  & breast,  $q = 2e-3$ ), SOX9 (colorectal & pan-cancer,  $q = 0$ ), TGIF1 (pan-cancer,  $q = 3e-3$  &  
165 colorectal  $q = 6e-3$ ), ZFP36L1 (urinary tract  $q = 3e-4$ , pan-cancer  $q = 9e-4$ ) and ZFP36L2 (colorectal &  
166 pan-cancer,  $q = 0$ ), HLA-B (lymphoid,  $q = 5e-5$ ), MGA (pan-cancer,  $q = 4e-03$ ), KMT2B (skin,  $q = 3e-3$ ) and  
167 RARG (urinary tract 8e-4). None of these genes are currently included in the COSMIC Cancer Gene  
168 Census<sup>23</sup>. ZFP36L1 and ZFP36L2 are of particular interest as these genes are zinc-finger proteins  
169 that normally play a repressive regulatory role in cell proliferation, presumably through a cyclin D  
170 dependent and p53 independent pathway<sup>24</sup>. ZFP36L2 is also independently found as a significantly  
171 deleted gene in our cohort, most prominently in colorectal and prostate cancers.

172 We also searched for genes that were significantly amplified or deleted (Supplementary Table  
173 4). CDKN2A and PTEN were the most significantly deleted genes overall, but many of the top genes  
174 involved common fragile sites (CFS) particularly FHIT and DMD, deleted in 5% and 4% of samples,  
175 respectively. The role of CFS in tumorigenesis is unclear and they are frequently treated as  
176 passenger mutations reflecting localized genomic instability<sup>25</sup>. However, the uneven distribution of the  
177 deletions across cancer types may indicate that some of these could be genuine tumor-type specific  
178 cancer drivers. For example, we find deletions in DMD to be highly enriched in esophageal tumors  
179 (38% deleted), GIST (Gastro-Intestinal Stromal Tumors; 24%) and pancreatic NET (41%), which is  
180 consistent with a recent study that indicated DMD as a TSG in cancers with myogenic programs<sup>26</sup>.  
181 We also identified several significantly deleted genes not reported previously, including MLLT4 (13  
182 samples) and PARD3 (9 samples).

183 Unlike homozygous deletions, amplification peaks tend to be broad and often encompass  
184 large number of genes, making identification of the amplification target challenging. However, SRY-  
185 related high-mobility group box 4 gene (SOX4, 6p22.3) stands out as a significantly amplified single  
186 gene peak (26 amplifications) and is highly enriched in urinary tract cancers (19% of samples highly  
187 amplified) (Extended Data Fig. 4). SOX4 is known to be over-expressed in prostate, hepatocellular,  
188 lung, bladder and medulloblastoma cancers with poor prognostic features and advanced disease  
189 status and is a modulator of the PI3K/Akt signaling<sup>27</sup>.

190 Also notable was a broad amplification peak of 10 genes around ZMIZ1 at 10q22.3 (32  
191 samples) which has not previously been reported. ZMIZ1 is a member of the Protein Inhibitor of  
192 Activated STAT (PIAS)-like family of coregulators and is a direct and selective cofactor of Notch1 in T-  
193 cell development and leukemia<sup>28</sup>. CDX2, previously identified as an amplified lineage-survival  
194 oncogene in colorectal cancer<sup>29</sup>, is also highly amplified in our cohort with 20 out of 22 amplified  
195 samples found in colorectal cancer, representing 5.4% of all colorectal samples.

196 **Driver mutation catalog**

197 We created a comprehensive catalog of all cancer driver mutations across all tumors in our  
198 cohort and all variant classes similar as described before<sup>30</sup> (N. Lopez, personal communication). To  
199 do this, we combined our SMG discovery efforts with those from Martincorena et al.<sup>21</sup> and a panel of  
200 well known cancer genes (Cosmic Curated Genes)<sup>31</sup>, and searched for fusions and TERT promoter  
201 mutations in our cohort. Accounting for the proportion of SNV and INDELS estimated by dNdScv to be  
202 passengers, we found 13,433 driver events among the 20,352 identified mutations in the combined  
203 gene panel ([Supplementary table 5](#)). This includes 7,423 coding mutation drivers, 615 non-coding  
204 point mutation drivers, 2,715 homozygous deletions (25% of which are in common fragile sites), 2,393  
205 focal amplifications and 287 fusion events. To facilitate analysis of variants of unknown significance  
206 (VUS) at a per patient level, we calculated a sample specific likelihood for each point mutation to be a  
207 driver taking into account the TMB of the sample as well as the biallelic inactivation status of the gene  
208 for TSG and hotspot positions in oncogenes. Of note, only 52% of point mutations in driver genes  
209 were predicted to be genuine driver events. Predictions of pathogenic variant overlap with known  
210 biology, e.g. clustering of benign missense variants in the 3' half of the APC gene ([Extended Data Fig.](#)  
211 [6b](#)) fits with the absence of FAP-causing germline variants in this part of the gene<sup>32</sup>.

212 Overall, the catalog matches previous inventories of cancer drivers. TP53 (52% of samples),  
213 CDKN2A (21%), APC (16%), PIK3CA (16%), KRAS (14%), PTEN (13%) and TERT (12%) were the  
214 most common driver genes together making up 25% of all the driver mutations in the catalog ([Fig. 3](#)).  
215 However, all of the ten most prevalent driver genes in our cohort were reported at a higher rate than  
216 for primary cancers<sup>33</sup>, which may reflect the more advanced disease state. AR and ESR1 were found  
217 as outliers with driver mutations in 44% of prostate and 18% of breast cancers, respectively. Both  
218 genes are linked to resistance to hormonal therapy, a common treatment for these tumor types, and  
219 have been previously reported as enriched in advanced metastatic cancer<sup>11</sup> but are identified at  
220 higher rates in this study.

221 Looking at a per patient level, the mean number of mutated driver genes per patient was 5.6  
222 (including any type of mutational event), with highest rate in urinary tract tumors (mean rate = 8.0) and  
223 the lowest in NET (mean rate = 2.8) ([Fig. 4](#)). Esophageal and stomach tumors also had highly  
224 elevated driver counts, largely due to a much higher rate of deletions in CFS genes (mean rate = 1.5  
225 for stomach, 1.7 for esophageal) compared to other cancer types (pan-cancer mean rate = 0.3).  
226 Fragile sites aside, the differential rates of drivers between cancer types in each variant class do  
227 correlate with the relative mutational load, with the exception of skin cancers, which have a lower than  
228 expected number of SNV drivers ([Extended Data Fig. 2](#)).

229 In 98.5% of all samples at least one driver mutation was found. Of the 36 samples in which no  
230 tumor driver mutation was identified, 18 were NET of the small intestine (representing 46.5% of all  
231 patients in this category). This likely indicates that small intestine NETs have a distinct set of drivers  
232 that are not captured yet in any of the cancer gene resources used and are also not prevalent enough  
233 in our relatively small NET cohort to be detected as significant.

234 The number of amplified driver genes varied significantly between cancer types with highly  
235 elevated rates per sample in breast cancer (mean = 2.0), esophageal, urinary tract and stomach (all  
236 mean = 1.7) cancers and nearly no amplification drivers in kidney cancer (mean = 0.1) and none in  
237 the mesothelioma cohort ([Extended Data Fig. 7a](#)). In tumor types with high rates of amplifications,  
238 these amplifications are generally found across a broad spectrum of oncogenes ([Extended Data Fig.](#)  
239 [7b](#)), suggesting there are mutagenic processes active in these tissues that favor amplifications, rather  
240 than tissue-specific selection of individual driver genes. AR and EGFR are notable exceptions, with  
241 highly selective amplifications in prostate, and in CNS and lung cancers, respectively, in line with  
242 previous reports<sup>17,34,35</sup>. Intriguingly, we also found two-fold more amplification drivers in samples with  
243 WGD ([Extended Data Fig. 7c](#)) despite amplifications being defined as relative to the average sample  
244 ploidy.

245 Analysis of known fusion drivers and promiscuous fusion partners identified 175 in-frame  
246 fusion in genic regions and 92 cis-activating fusions involving repositioning of regulatory elements,  
247 and 20 in-frame intragenic deletions where one or more exons were deleted. ERG (89 samples),

248 BRAF(17 samples), ERBB4 (16 samples), ALK(12 samples), NRG1(9 samples) and ETV4 (7  
249 samples) were the most commonly observed 3' partners together making up more than half of the  
250 fusions. 77 of the 89 ERG fusions were TMPRSS2-ERG affecting 37% of all prostate cancer samples  
251 in the cohort. 153 fusion pairs were not previously recorded in CGI, OncoKb, COSMIC or CIViC<sup>31,36-38</sup>.  
252 A novel recurrent KMT2A-BCOR fusion was observed in 2 samples (sarcoma and stomach cancer)  
253 and there were also 5 recurrent localized fusions resulting from adjacent gene pairs: NSD1-ZNF346  
254 (3 samples), FGFR2-ATE1 (2 samples), KMT2A-BCOR(2 samples), AGM1-ETV1 (2 samples), and  
255 BCR-GNAZ (2 samples).

256 Only promoter mutations in TERT were included in the study due to the current lack of robust  
257 evidence for other recurrent oncogenic non-coding mutations<sup>39</sup>. A total of 257 variants were found at 5  
258 known recurrent variant hotspots<sup>11</sup>, and included in the driver catalog.

## 259 **Oncogene hotspots and novel activating variants**

260 We found that the 70% of driver mutations in oncogenes occur at or within 5 nucleotides of  
261 already known pathogenic mutational hotspots ([Extended Data Fig. 6a](#)). In the six most prevalent  
262 oncogenes (KRAS, PIK3CA, BRAF, NRAS, CTNNB1 & ESR1) the rate was 96% ([Fig. 5](#)).  
263 Furthermore, in many of the key oncogenes, we also found several likely activating but non-canonical  
264 variants near known mutational hotspots ([Fig. 5](#)). For example, we found activating MNVs in the well  
265 known BRAF V600 hotspot (22 cases), but also novel non-hotspot MNVs in KRAS (8 cases) and  
266 NRAS (4 cases) ([Extended Data Fig 6b](#)).

267 In-frame indels were even more striking, since despite being exceptionally rare overall (mean  
268 = 1.7 per sample), we found an excess in known oncogenes including PIK3CA (19 cases), ERBB2  
269 (10 cases) and BRAF(8 cases) frequently occurring at or near known hotspots. Notably, many of the  
270 in-frame indels are enriched in specific tumor types. For instance, all 18 KIT in-frame indels were  
271 found in sarcomas, 6 out of 8 MUC6 in-frame indels in esophageal tumors, and 6 of 10 ERBB2 in-  
272 frame indels in lung tumors. Finally, we identified 10 in-frame indels in FOXA1, which are highly  
273 enriched in prostate cancer (7 of 10 cases) and clustered in two locations that were not previously  
274 associated with pathogenic mutations<sup>40</sup>.

275 In CTNNB1 we identified an interesting novel recurrent in-frame deletion of the complete exon  
276 3 in 12 samples, 9 of which are colorectal cancers. Surprisingly, these deletions were homozygous  
277 but suspected to be activating as CTNNB1 normally acts as an oncogene in the WNT/beta-catenin  
278 pathway and none of these nine colorectal samples had any APC driver mutations.

## 279 **Biallelic tumor suppressor gene inactivation**

280 Our results strongly support the Knudson two-hit hypothesis<sup>41</sup> for tumor suppressor genes  
281 with 80% of all TSG drivers explained by biallelic inactivation (i.e. either by homozygous deletion  
282 (32%), multiple somatic point mutations (7%), or a point mutation in combination with LOH (41%)).  
283 This rate is the highest observed in any large-scale cancer WGS study. For many key tumor  
284 suppressor genes the biallelic inactivation rate is almost 100% (more specifically: TP53 (93%),  
285 CDKN2A (97%), RB1 (94%), PTEN (92%) and SMAD4 (96%); [Fig. 3b](#)), suggesting that biallelic  
286 inactivation of these genes is a strict requirement for metastatic cancer.

287 Other prominent TSGs, however, have lower biallelic rates, including ARID1A (55%), KMT2C  
288 (49%) and ATM (49%). It is unclear whether we systematically missed the second hit in these cases,  
289 as this could potentially be mediated through non-mutational epigenetic mechanisms<sup>42</sup> or unknown  
290 pathogenic germline variants, or if these genes impact on tumorigenesis by a haploinsufficiency  
291 mechanism<sup>43</sup>.

## 292 **Clonal and subclonal variants**

293 To obtain insight into ongoing tumor evolution dynamics, we examined the clonality of all  
294 variants. Surprisingly, only 6.5% of all SNV, MNV & INDELS across our cohort and just 3.7% of our

295 driver point mutations were found to be subclonal (Fig. 6). The low proportion of samples with  
296 subclonal variants could be partially due to the detection limits of the sequencing approach,  
297 particularly for low purity samples, even with our relatively high WGS sequencing depth. However,  
298 even for samples with purities higher than 80% the proportion of subclonal variants only reaches  
299 10.2% (Fig. 6c). This relatively high degree of tumor homogeneity may be in part attributed to the fact  
300 that nearly all biopsies were obtained by a core needle biopsy, which results in highly localized  
301 sampling, but is nevertheless much lower compared to previous observations in primary cancers<sup>2</sup>.

302 In the 111 patients with independently collected repeat biopsies from the same patient  
303 (Supplementary Table 6) we found 11% of all SNVs to be subclonal. Whilst 76% of clonal variants  
304 were shared between biopsies, less than 30% of the subclonal variants were shared. Together, the  
305 low rate of subclonal variants, and the observation that a very high proportion are private to a local  
306 metastasis, suggest a model where individual metastatic lesions are dominated by a single clone at  
307 any one point in time and that more limited tumor evolution and subclonal selection takes places after  
308 distant metastatic seeding.

### 309 **Co-occurrence of Drivers**

310 We examined the pairwise co-occurrence of driver gene mutations per cancer type and found  
311 ten combinations of genes that were significantly mutually exclusively mutated, and ten combinations  
312 of genes which were significantly co-occurringly mutated (excluding pairs of genes on the same  
313 chromosome which are frequently co-amplified or co-deleted) (Fig. 7). The 20 significant findings  
314 include previously reported co-occurrence of mutated DAX|MEN1 in pancreatic NET ( $q=0.0007$ ), and  
315 CDH1|SPOP in prostate tumors ( $q=0.0005$ ), as well as negative associations of mutated genes within  
316 the same signal transduction pathway such as KRAS|BRAF ( $q=4e-4$ ) & KRAS|NRAS ( $q=0.009$ ) in  
317 colorectal cancer, BRAF|NRAS in skin cancer ( $q=6e-12$ ), CDKN2A|RB1 in lung cancer ( $q=8e-5$ ) and  
318 APC|CTNNB1 in colorectal cancer ( $q=8e-6$ ). APC is also strongly negatively correlated with BRAF  
319 ( $q=1e-4$ ) and RNF43 ( $q=2e-5$ ) which together are characteristic of the serrated molecular subtype of  
320 colorectal cancers<sup>44</sup>. We also found that SMAD2|SMAD3 are highly positively correlated in colorectal  
321 cancer ( $q=0.02$ ), mirroring a result reported previously in a large cohort of colorectal cancers<sup>45</sup>.

322 In breast cancer, we found a number of significant novel relationships, including a positive  
323 relationship for GATA3|VMP1( $q=1e-4$ ) and FOXA1|PIK3CA ( $q=2e-3$ ), and negative relationships for  
324 ESR1|TP53 ( $q=9e-4$ ) and GATA3|TP53 ( $q=2e-3$ ).

### 325 **Actionability**

326 We analyzed opportunities for biomarker-based treatment for all patients by mapping driver  
327 events to three clinical annotation databases: CGI<sup>38</sup>, CIViC<sup>36</sup> and OncoKB<sup>37</sup>. In 1,485 patients (62%)  
328 at least one ‘actionable’ event was identified (Supplementary Table 7). While these numbers are in  
329 line with results from primary tumors<sup>30</sup>, longitudinal studies will be required to conclude if genomic  
330 analyses for therapeutic guidance should be repeated when a patient experiences progressive  
331 disease. Half of the patients with an actionable event (31% of total) contained a biomarker with a  
332 predicted sensitivity to a drug at level A (approved anti-cancer drugs) and lacked any known  
333 resistance biomarkers for the same drug (Fig. 8). In 13% of patients the suggested therapy was a  
334 registered indication, while in 18% of cases it was outside the labeled indication. In a further 31% of  
335 patients a level B (experimental therapy) biomarker was identified. The actionable biomarkers  
336 spanned all variant classes including 1,815 SNVs, 48 MNVs, 195 indels, 745 CNAs, 68 fusion genes  
337 and 60 patients with microsatellite instability.

338 Tumor mutation burden is an important emerging biomarker for response to immune  
339 checkpoint inhibitor therapy<sup>46</sup> as it is a proxy for the amount of neo-antigens in the tumor cells. For  
340 NSCLC it has been shown in at least 2 subgroup analyses of large phase III trials that both PFS and  
341 OS are significantly improved with first line immunotherapy as compared to chemotherapy for patients  
342 whose tumors have a TMB >10 mutations per Mb<sup>47,48</sup>. Although various clinical studies based on this  
343 parameter are currently emerging, TMB was not yet included in the above actionability analysis.

344 However, when applying the same cut-off to all samples in our cohort, an overall 18% of patients  
345 would qualify, varying from 0% for liver, mesothelioma and ovarian cancer patients to more than 50%  
346 of lung and skin cancer patients ([Extended Data Fig. 3b](#)).

347 **Discussion**

348 Genomic testing of tumors faces numerous challenges in meeting clinical needs, including i)  
349 the interpretation of variants of unknown significance (VUS), ii) the steadily expanding universe of  
350 actionable genes, often with an increasingly small fraction of patients affected (e.g. NRG-1<sup>49</sup> and  
351 NTRK fusions<sup>50</sup> in less than 2% of all patients), and iii) the development of advanced genome-derived  
352 biomarkers such as tumor mutational load, DNA repair status and mutational signatures. Our results  
353 demonstrate in several ways that WGS analyses of metastatic cancer can provide novel and relevant  
354 insights and be instrumental in addressing some of these key challenges in cancer precision  
355 medicine.

356 First, our systematic and large-scale pan-cancer analyses on metastatic cancer tissue  
357 allowed for the identification of several novel (cancer type-specific) cancer drivers and mutation  
358 hotspots. Second, the driver catalog analyses can be used to alleviate the problem of VUS  
359 interpretation<sup>30</sup> by leveraging previously identified pathogenic mutations (accounting for more than  
360 2/3rds of oncogenic point-mutation drivers) and by careful analysis of the biallelic inactivation of  
361 putative TSG which accounts for over 80% of TSG drivers in metastatic cancer. Third, we  
362 demonstrate the importance of accounting for all types of variants, including large scale genomic  
363 rearrangements (via fusions and copy number alteration events), which account for more than half of  
364 all drivers, but also activating MNV and INDELs which we have shown are commonly found in many  
365 key oncogenes. Fourth, we have shown that using WGS, even with very strict variant calling criteria,  
366 we could find driver variants in more than 98% of all metastatic tumors, including putatively actionable  
367 events in a clinical and experimental setting for up to 62% of patients.

368 Although we did not find metastatic tumor genomes to be fundamentally different to primary  
369 tumors in terms of the mutational landscape or genes driving advanced tumorigenesis, we described  
370 characteristics that could contribute to therapy responsiveness or resistance in individual patients. In  
371 particular we showed that WGD is a more pervasive element of tumorigenesis than previously  
372 understood affecting over half of all metastatic cancers. We also found metastatic lesions to be less  
373 heterogeneous than primary tumors, supporting a model for metastasis of clonal seeding followed by  
374 rapid expansion.

375 It should be noted that differences between WGS cohorts should be interpreted with some  
376 caution as inevitable differences between experimental and computational approaches may impact on  
377 observations and can only be addressed in longitudinal studies including the different stages of  
378 disease. Furthermore, the HMF cohort includes a mix of treatment-naive metastatic patients and  
379 patients who have undergone (extensive) prior systemic treatments. While this may impact on specific  
380 tumor characteristics, it also provides opportunities for studying treatment response and resistance as  
381 this data is recorded in the studies.

382 Finally, we believe the resource described here is a valuable complementary resource to  
383 comparable whole genome sequencing-based resources of primary tumors in advancing fundamental  
384 and translational cancer research. Therefore, all non-privacy sensitive data is publicly available  
385 through a local interface developed by ICGC<sup>51</sup> (work in progress) and all other data is made freely  
386 available for scientific research by a controlled access mechanism (see  
387 [www.hartwigmedicalfoundation.nl](http://www.hartwigmedicalfoundation.nl) for details).

388

389 **Acknowledgements**

390 We thank the Hartwig Foundation and Barcode for Life for financial support of clinical studies and  
391 WGS analyses. Development of the data portal was supported by a grant from KWF  
392 Kankerbestrijding (HMF2017-8225, GENONCO). We are particularly grateful to all patients, nurses  
393 and medical specialists for their essential contributions making this study possible. We would like to  
394 specifically thank Hans van Snellenberg for operational management of the Hartwig Medical  
395 Foundation. We would like to thank Stefan Willems, Wendy de Leng, Alexander Hoischen and  
396 Winand Dinjens for support with pathology assessments and mutation validations and Jeroen de  
397 Ridder, Wigard Kloosterman and Harmen van de Werken for critically reading the manuscript.

398

399

400 **Figure Legends**

401 **Figure 1: Mutational load of metastatic cancer per tumor type**

402 a) The number of samples of each tumor type cohort. Tumor types are ranked from lowest to highest  
403 overall mutation burden (TMB)  
404 b) Violin plot showing age distribution of each tumor type with 25th, 50th and 75th percentiles marked.  
405 c)-d) cumulative distribution function plot of mutational load for each tumor type for SNV and MNV (c)  
406 and INDEL and SV (d). The median for each cohort is indicated with a vertical line.  
407 e)-h) Mutational context or variant subtype per individual sample for each of (e) Single Nucleotide  
408 Variant (SNV), (f) Multi Nucleotide Variant (MNV), (g) Insertion/Deletion (INDEL), (h) Structural  
409 Variant (SV). Each column chart is ranked within tumor type by mutational load in that variant class.  
410 MNVs are classified by the dinucleotide substitution with NN referring to any dinucleotide  
411 combination. SVs are classified by type: INV = inversion, DEL = deletion, DUP = tandem duplication,  
412 TRL = translocation, INS = insertion.

413

414 **Figure 2: Copy number landscape of metastatic cancer**

415 a)-c) Proportion of samples with amplification and deletion events by genomic position per cohort -  
416 pan-cancer (a), kidney (b) and central nervous system (CNS) (c). Inner ring shows % of tumors with  
417 homozygous deletion (orange), LOH and significant loss (copy number < 0.6x sample ploidy - dark  
418 blue) and near copy neutral LOH (light blue). Outer ring shows % of tumors with high level  
419 amplification (>3x sample ploidy - orange), moderate amplification (>2x sample ploidy - dark green)  
420 and low level amplification (>1.4x amplification - light green). Scale on both rings is 0-100% and  
421 inverted for the inner ring. The most frequently observed high-level gene amplifications (black text)  
422 and homozygous deletions (red text) are shown.  
423 d) Proportion of tumors with a whole genome duplication event (dark blue) grouped by tumor type.  
424 e) Average sample ploidy distribution over the complete cohort. Samples with a WGD event (true) are  
425 shown in darker blue.

426

427 **Figure 3: Most prevalent driver genes in metastatic cancer**

428 Thirty most prevalent mutated oncogenes (a) and tumor suppressor genes (TSG) (b). From left to  
429 right, the heatmap shows the % of samples in each cancer type which are found to have each gene  
430 mutated; absolute bar chart shows the pan-cancer % of samples with the given gene mutated; relative  
431 bar chart shows the breakdown by type of alteration. For TSG only the % of samples with a driver in  
432 which the gene is found biallelically inactivated is also shown.

433

434 **Figure 4: Drivers per sample by tumor type**

435 a) Violin plot showing the distribution of number of drivers per sample grouped by tumor type. Black  
436 dot indicates the mean value.  
437 b) Relative bar chart showing the breakdown per cancer type of the type of alteration.

438

439 **Figure 5: Oncogenic Hotspots**

440 Count of driver point mutations by variant type. Known pathogenic mutations curated from external  
441 databases are categorized as hotspot mutations. Mutations within 5 bases of a known pathogenic  
442 mutation are shown as near hotspot and all other mutations are shown as non-hotspot.

443

444 **Figure 6: Subclonality**

445 a) Count of samples per tumor purity bucket  
446 b) Violin plot showing the percentage of point mutations which are subclonal in each purity bucket.  
447 Black dot indicates the mean.  
448 c) Percentage of driver point mutations that are subclonal in each purity bucket.

449

450 **Figure 7: Driver co-occurrence**

451 a) Driver pairs which are significantly positively (on the right) or negatively (on the left) correlated in  
452 individual tumor types sorted by q-value. Color indicates tumor type as depicted below the chart.

453

454 **Figure 8: Actionability**

455 Percentage of samples in each cancer type with an actionable mutation based on data in CGI, CIViC  
456 and OncoKB knowledgebases. Level 'A' represents presence of biomarkers with either an approved  
457 therapy or guidelines and level B represents biomarkers having strong biological evidence or clinical  
458 trials indicating they are actionable. On label indicates treatment registered by federal authorities for  
459 that tumor type, while off-label indicates a registration for other tumor types.

460

461 **Extended Data Figures and Tables**

462

463 **Extended Data Figure 1: Hartwig sample workflow, biopsy locations and sequence coverage**

464 A) Sample workflow from patient to high-quality WGS data. A total of 4,018 patients were enrolled in  
465 the study between April 2016 and April 2018. For 9% of patients no blood and/or biopsy material was  
466 obtained, mostly because conditions of patients prohibited further study participation. Up to 4 fresh-  
467 frozen biopsies per patient were received, which were sequentially analyzed to identify a biopsy with  
468 more than 30% tumor cellularity as determined by routine histology assessment. For 859 patients no  
469 suitable biopsy was obtained and 2,796 patients were further processed for WGS. 44 and 29 samples  
470 failed in either DNA isolation or library preparation and raw WGS data quality QC, respectively. For an  
471 additional 385 samples the WGS data was of good quality, but the tumor purity determination based  
472 on WGS data (PURPLE) was less than 20% making reliable and comprehensive somatic variant  
473 calling and were therefore excluded. Eventually, 2,338 tumor-normal sample pairs with high-quality  
474 WGS data were obtained, which were supplemented with 182 pairs from pre-April 2016, adding up to  
475 2,520 tumor normal pairs that were included in this study.

476 B) Breakdown of cohort by biopsy location. Tumor biopsies were taken from a broad range of  
477 locations. Primary tumor type is shown on the left and the biopsy location on the right.

478 c) Distribution of sample sequencing depth for tumor and blood reference.

479

480 **Extended Data Figure 2: Impact of downsampling on variant calling**

481 Comparison of variant calling of 10 randomly selected samples at normal depth and 50%  
482 downsampled for purity (a), SNV counts (b), SV counts (c), Ploidy (d), MNV counts (e) and INDEL

483 counts (f). For the panels B, C, E and F, the black dots represent the % reduction per sample of  
484 counts (right axis) and the dotted line represents the average % reduction across all tested samples.  
485

486 **Extended Data Figure 3: Mutational load, genome wide analyses and drivers**

487 a) Proportion of samples by cancer type classified as microsatellite unstable (MSI)  
488 b) Proportion of samples with a high mutational burden (TMB > 10 SNV / Mb)  
489 c) Scatter plot of INDEL vs SNV mutational load. MSI and 'high TMB' thresholds are indicated.  
490 d-f) Mutational load vs driver rate for SNV (d), INDEL (e) and SV (f) grouped by cancer type. MSI  
491 samples were excluded.

492

493 **Extended Data Figure 4: Copy Number profile per cancer types**

494 Proportion of samples with amplification and deletion events by genomic position per cancer type.  
495 Inner ring shows % of tumors with homozygous deletion (red), LOH and significant loss (copy number  
496 < 0.6x sample ploidy - dark blue) and near copy neutral LOH (light blue). Outer ring shows % of  
497 tumors with high level amplification (>3x sample ploidy - orange), moderate amplification (>2x sample  
498 ploidy - dark green) and low level amplification (>1.4x amplification - light green). Scale on both rings  
499 are 0-100%, inverted for inner ring. The most frequently observed high level gene amplifications  
500 (black text) and homozygous deletions (red text) are labelled.

501

502 **Extended Data Figure 5: Significantly mutated genes**

503 Tile chart showing genes found to be significantly mutated per cancer type cohort and pan-cancer  
504 using dNdScv. Gene names marked in orange are also significant in Martincorena et al<sup>21</sup>, but not  
505 found in COSMIC curated or census. Gene names marked in red are novel in this study.

506

507 **Extended Data Figure 6: Coding mutation profiles by driver gene**

508 Location and driver classification of all coding mutations (SNVs and indels) in oncogenes (a) and  
509 tumor suppressor genes (TSG) (b) in the driver catalog. The lollipops on the chart show the location  
510 (coding sequence coordinates) and count of mutations for all candidate drivers. Dotted lines show the  
511 count of all variants at each location while solid lines represent the sum of the driver likelihoods  
512 assigned to that specific variant. TSG variants are shaded by their biallelic inactivation status. 'Mixed'  
513 means that at least 2 samples with the variant have different biallelic statuses. The bars on the right of  
514 the chart show the estimated number of drivers (calculated as the sum of driver likelihoods) and  
515 passenger variants in the gene by cancer type.

516

517 **Extended Data Figure 7: Amplifications**

518 a) Mean rate of amplification drivers per cancer type  
519 b) Number of amplification drivers per gene showing the breakdown by cancer type  
520 c) Mean rate of drivers per variant type for samples with and without WGD.

521

522 **Supplementary Table 1:** Overview of contributing organizations and local principal investigators.

523

524 **Supplementary Table 2:** Overview of cohort and sample characteristics

525

526 **Supplementary Table 3:** Pan-cancer and cancer type-specific dNdScv results

527

528 **Supplementary Table 4:** Recurring amplifications (a) and deletions (b) and associated target genes

529

530 **Supplementary Table 5:** Driver catalog

531

532 **Supplementary Table 6:** Overview of patients with multiple biopsies

533

534 **Supplementary Table 7:** Actionable mutations

535 References

- 536 1. Klein, C. A. Selection and adaptation during metastatic cancer progression. *Nature* **501**, 365–372  
537 (2013).

538 2. McGranahan, N. & Swanton, C. Clonal Heterogeneity and Tumor Evolution: Past, Present, and  
539 the Future. *Cell* **168**, 613–628 (2017).

540 3. Cancer Genome Atlas Research, Network *et al.* The Cancer Genome Atlas Pan-Cancer analysis  
541 project. *Nat. Genet.* **45**, 1113–1120 (2013).

542 4. International Cancer Genome, Consortium *et al.* International network of cancer genome  
543 projects. *Nature* **464**, 993–998 (2010).

544 5. Grobner, S. N. *et al.* The landscape of genomic alterations across childhood cancers. *Nature*  
545 **2018/03/01**, (2018).

546 6. Ma, X. *et al.* Pan-cancer genome and transcriptome analyses of 1,699 paediatric leukaemias and  
547 solid tumours. *Nature* **2018/03/01**, (2018).

548 7. Hyman, D. M., Taylor, B. S. & Baselga, J. Implementing Genome-Driven Oncology. *Cell* **168**,  
549 584–599 (2017).

550 8. Yates, L. R. *et al.* Genomic Evolution of Breast Cancer Metastasis and Relapse. *Cancer Cell* **32**,  
551 169–184 e7 (2017).

552 9. Naxerova, K. *et al.* Origins of lymphatic and distant metastases in human colorectal cancer.  
553 *Science* **357**, 55–60 (2017).

554 10. Gundem, G. *et al.* The evolutionary history of lethal metastatic prostate cancer. *Nature* **520**, 353–  
555 357 (2015).

556 11. Zehir, A. *et al.* Mutational landscape of metastatic cancer revealed from prospective clinical  
557 sequencing of 10,000 patients. *Nat. Med.* **23**, 703–713 (2017).

558 12. Robinson, D. R. *et al.* Integrative clinical genomics of metastatic cancer. *Nature* **548**, 297–303  
559 (2017).

560 13. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–  
561 421 (2013).

562 14. Gryfe, R. *et al.* Tumor microsatellite instability and clinical outcome in young patients with  
563 colorectal cancer. *N. Engl. J. Med.* **342**, 69–77 (2000).

564 15. Taylor, A. M. *et al.* Genomic and Functional Approaches to Understanding Cancer Aneuploidy.  
565 *Cancer Cell* **33**, 676–689.e3 (2018).

566 16. Sato, Y. *et al.* Integrated molecular analysis of clear-cell renal cell carcinoma. *Nat. Genet.* **45**,  
567 860–867 (2013).

568 17. Brennan, C. W. *et al.* The somatic genomic landscape of glioblastoma. *Cell* **155**, 462–477 (2013).

569 18. Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat. Genet.* **45**, 1134–  
570 1140 (2013).

571 19. Carter, S. L. *et al.* Absolute quantification of somatic DNA alterations in human cancer. *Nat.*  
572 *Biotechnol.* **30**, 413–421 (2012).

573 20. Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of advanced cancers.  
574 *Nat. Genet.* (2018). doi:10.1038/s41588-018-0165-1

575 21. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**,  
576 1029–1041 e21 (2017).

577 22. Marusiak, A. A. *et al.* Recurrent MLK4 Loss-of-Function Mutations Suppress JNK Signaling to  
578 Promote Colon Tumorigenesis. *Cancer Res.* **76**, 724–735 (2016).

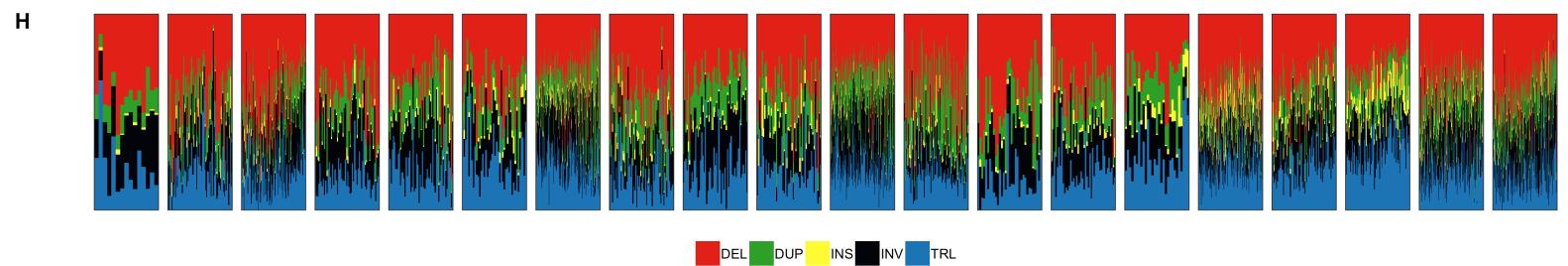
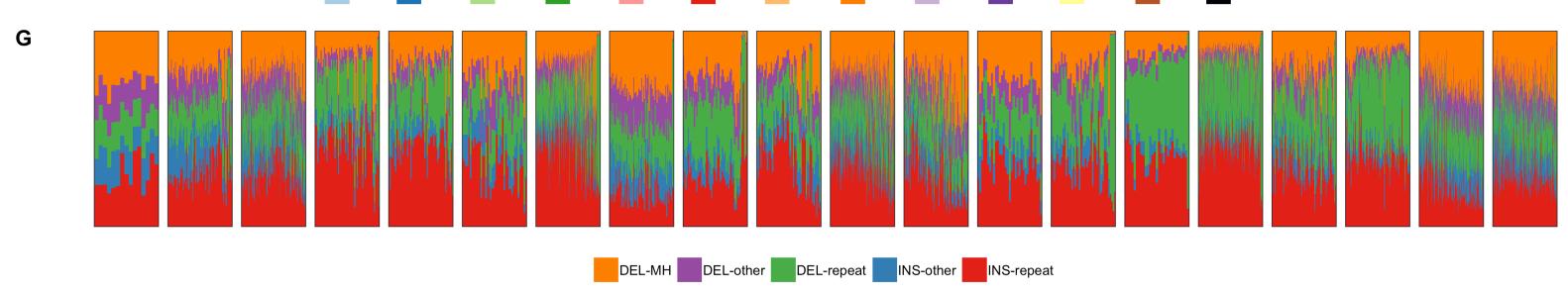
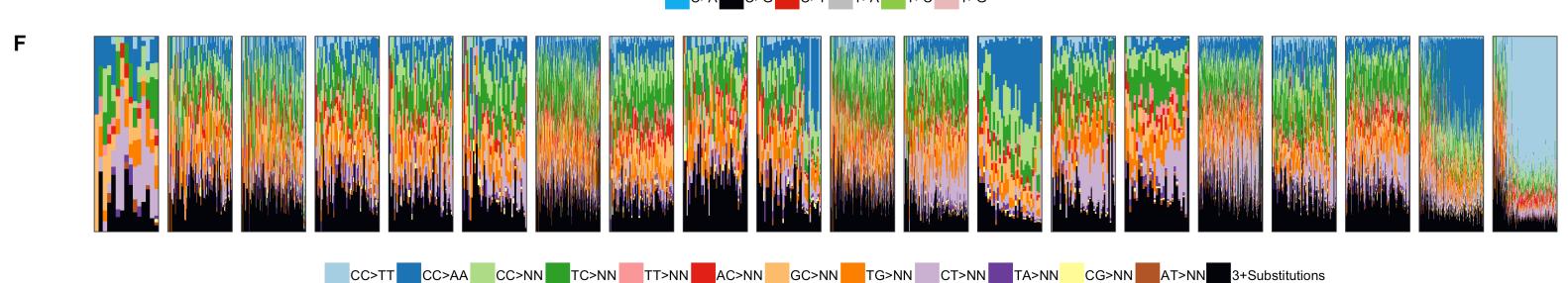
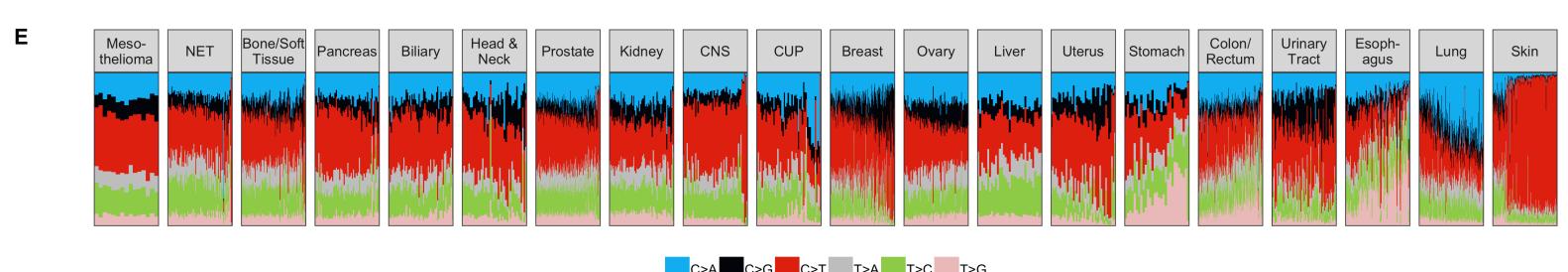
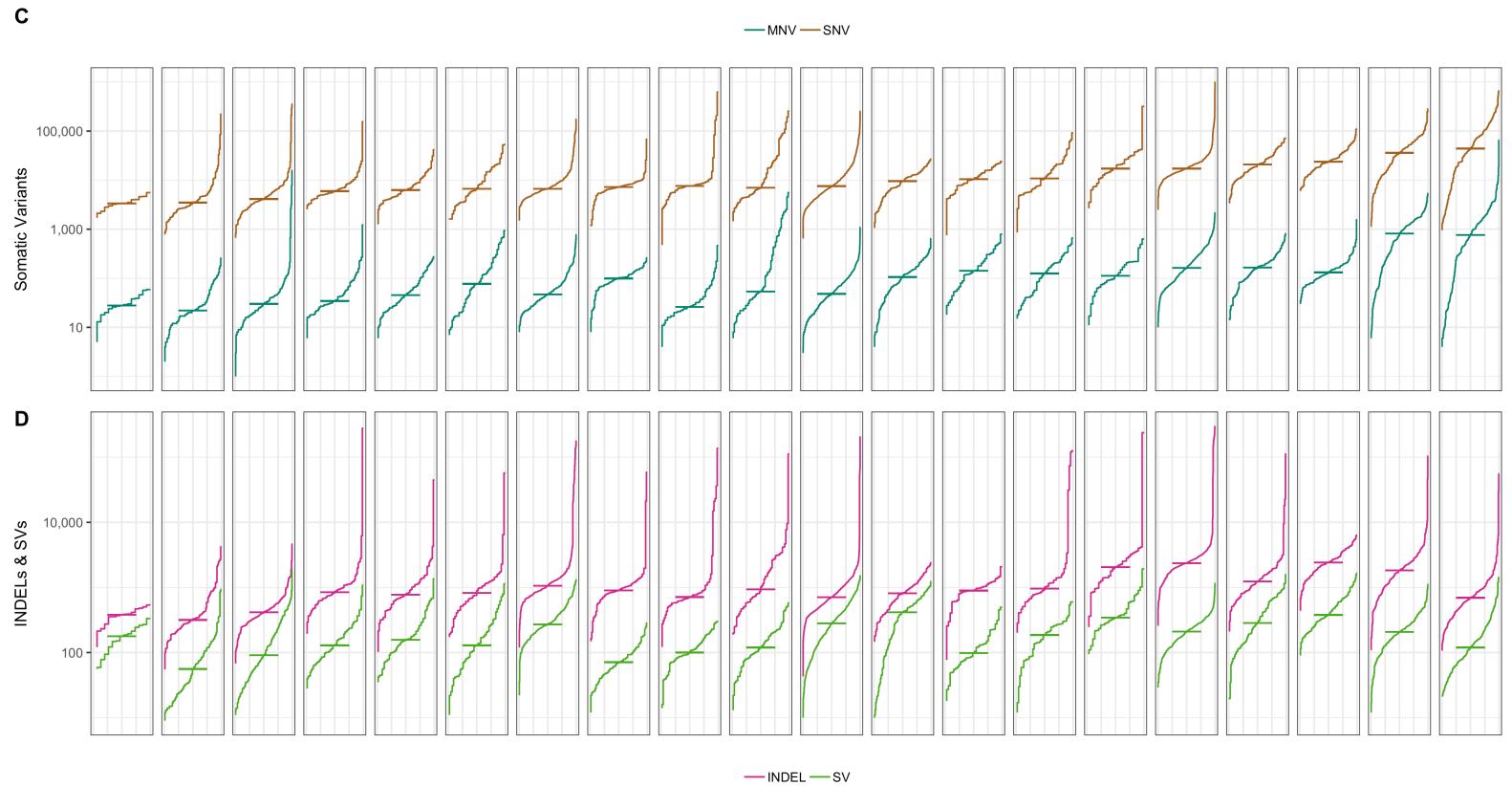
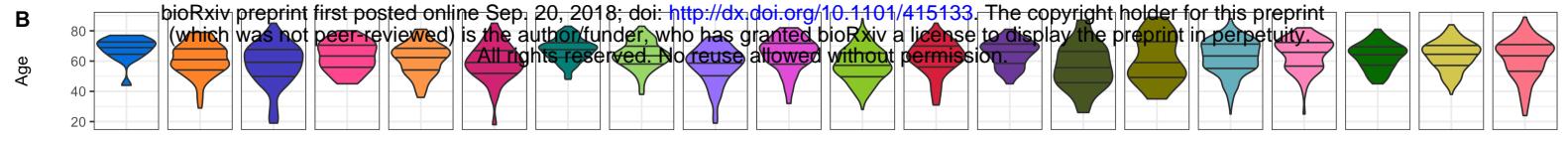
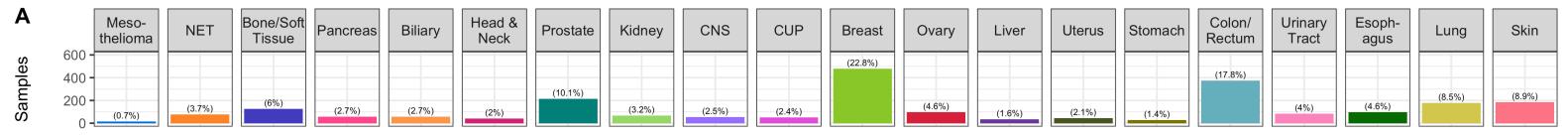
579 23. Forbes, S. A. *et al.* COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids Res.* **45**,  
580 D777–D783 (2017).

581 24. Suk, F.-M. *et al.* ZFP36L1 and ZFP36L2 inhibit cell proliferation in a cyclin D-dependent and p53-  
582 independent manner. *Sci. Rep.* **8**, 2742 (2018).

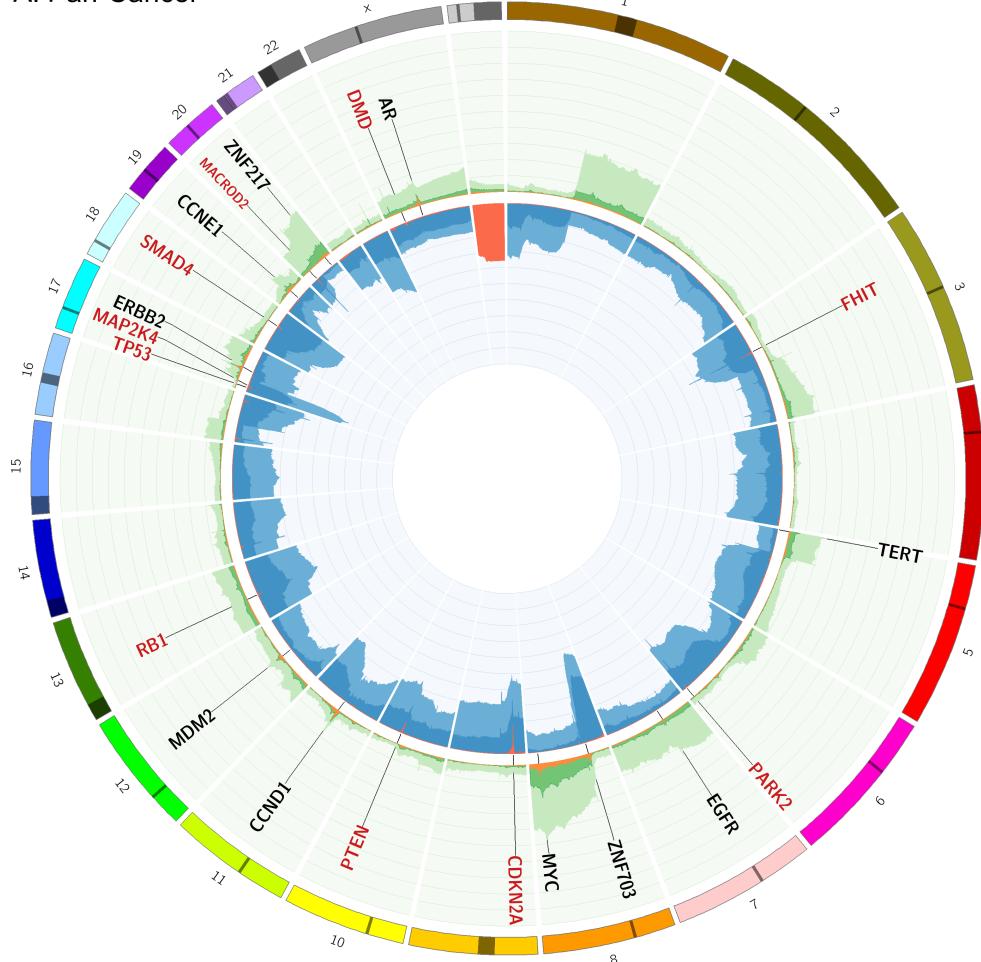
583 25. Glover, T. W., Wilson, T. E. & Arlt, M. F. Fragile sites in cancer: more than meets the eye. *Nat.*  
584 *Rev. Cancer* **17**, 489–501 (2017).

585 26. Wang, Y. *et al.* Dystrophin is a tumor suppressor in human cancers with myogenic programs.  
586 *Nat. Genet.* **46**, 601–606 (2014).

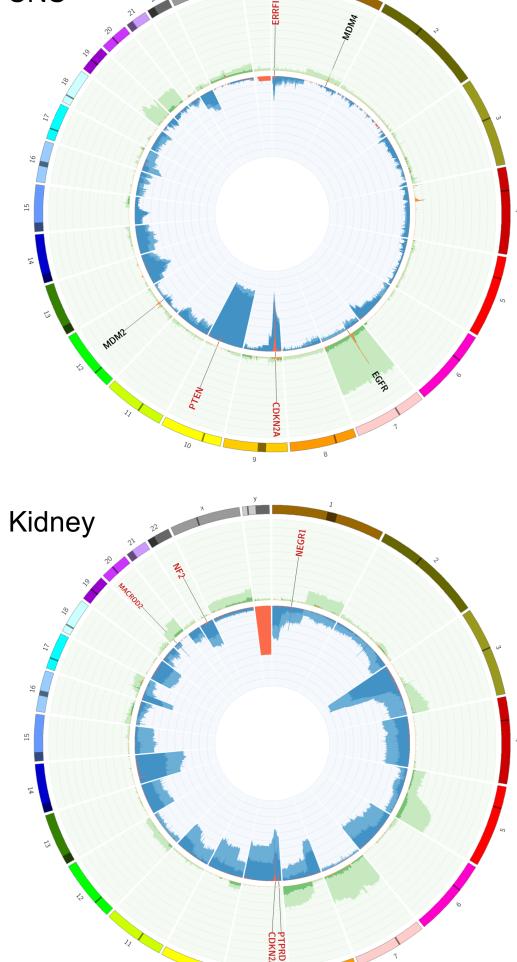
- 587 27. Mehta, G. A. *et al.* Amplification of SOX4 promotes PI3K/Akt signaling in human breast cancer. *Breast Cancer Res. Treat.* **162**, 439–450 (2017).
- 588 28. Pinnell, N. *et al.* The PIAS-like Coactivator Zmiz1 Is a Direct and Selective Cofactor of Notch1 in  
589 T Cell Development and Leukemia. *Immunity* **43**, 870–883 (2015).
- 590 29. Salari, K. *et al.* CDX2 is an amplified lineage-survival oncogene in colorectal cancer. *Proc. Natl.  
591 Acad. Sci. U. S. A.* **109**, E3196–205 (2012).
- 592 30. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *BioArchive* (2017).  
593 doi:10.1101/190330
- 594 31. Futreal, P. A. *et al.* A census of human cancer genes. *Nat. Rev. Cancer* **4**, 177–183 (2004).
- 595 32. Friedl, W. *et al.* Can APC mutation analysis contribute to therapeutic decisions in familial  
596 adenomatous polyposis? Experience from 680 FAP families. *Gut* **48**, 515–521 (2001).
- 597 33. Bailey, M. H. *et al.* Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell*  
598 **173**, 371–385.e18 (2018).
- 599 34. Viswanathan, S. R. *et al.* Structural Alterations Driving Castration-Resistant Prostate Cancer  
600 Revealed by Linked-Read Genome Sequencing. *Cell* (2018). doi:10.1016/j.cell.2018.05.036
- 601 35. Cancer Genome Atlas Research Network. Comprehensive genomic characterization of  
602 squamous cell lung cancers. *Nature* **489**, 519–525 (2012).
- 603 36. Griffith, M. *et al.* CIViC is a community knowledgebase for expert crowdsourcing the clinical  
604 interpretation of variants in cancer. *Nat. Genet.* **49**, 170–174 (2017).
- 605 37. Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol* **2017**,  
606 (2017).
- 607 38. Tamborero, D. *et al.* Cancer Genome Interpreter annotates the biological and clinical relevance  
608 of tumor alterations. *Genome Med.* **10**, 25 (2018).
- 609 39. Cuykendall, T. N., Rubin, M. A. & Khurana, E. Non-coding genetic variation in cancer. *Current  
610 Opinion in Systems Biology* **1**, 9–15 (2017).
- 611 40. Yang, Y. A. & Yu, J. Current perspectives on FOXA1 regulation of androgen receptor signaling  
612 and prostate cancer. *Genes Dis* **2**, 144–151 (2015).
- 613 41. Knudson, A. G., Jr. Mutation and cancer: statistical study of retinoblastoma. *Proc. Natl. Acad. Sci.  
614 U. S. A.* **68**, 820–823 (1971).
- 615 42. Schlicker, A., Michaut, M., Rahman, R. & Wessels, L. F. A. OncoScape: Exploring the cancer  
616 aberration landscape by genomic data fusion. *Sci. Rep.* **6**, 28103 (2016).
- 617 43. Davoli, T. *et al.* Cumulative haploinsufficiency and triplosensitivity drive aneuploidy patterns and  
618 shape the cancer genome. *Cell* **155**, 948–962 (2013).
- 619 44. Bond, C. E. *et al.* RNF43 and ZNRF3 are commonly altered in serrated pathway colorectal  
620 tumorigenesis. *Oncotarget* **7**, 70589–70600 (2016).
- 621 45. Fleming, N. I. *et al.* SMAD2, SMAD3 and SMAD4 mutations in colorectal cancer. *Cancer Res.* **73**,  
622 725–735 (2013).
- 623 46. Goodman, A. M. *et al.* Tumor Mutational Burden as an Independent Predictor of Response to  
624 Immunotherapy in Diverse Cancers. *Mol. Cancer Ther.* **16**, 2598–2608 (2017).
- 625 47. Hellmann, M. D. *et al.* Nivolumab plus Ipilimumab in Lung Cancer with a High Tumor Mutational  
626 Burden. *N. Engl. J. Med.* **378**, 2093–2104 (2018).
- 627 48. Carbone, D. P. *et al.* First-Line Nivolumab in Stage IV or Recurrent Non-Small-Cell Lung Cancer.  
628 *N. Engl. J. Med.* **376**, 2415–2426 (2017).
- 629 49. Fernandez-Cuesta, L. & Thomas, R. K. Molecular Pathways: Targeting NRG1 Fusions in Lung  
630 Cancer. *Clin. Cancer Res.* **21**, 1989–1994 (2015).
- 631 50. Laetsch, T. W. *et al.* Larotrectinib for paediatric solid tumours harbouring NTRK gene fusions:  
632 phase 1 results from a multicentre, open-label, phase 1/2 study. *Lancet Oncol.* **19**, 705–714  
633 (2018).
- 634 51. Zhang, J. *et al.* International Cancer Genome Consortium Data Portal--a one-stop shop for  
635 cancer genomics data. *Database* **2011**, bar026 (2011).
- 636



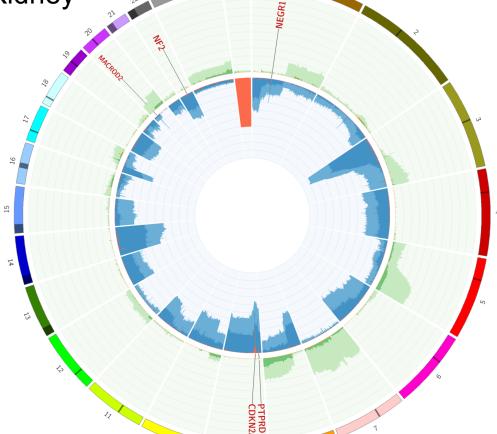
### A. Pan-Cancer



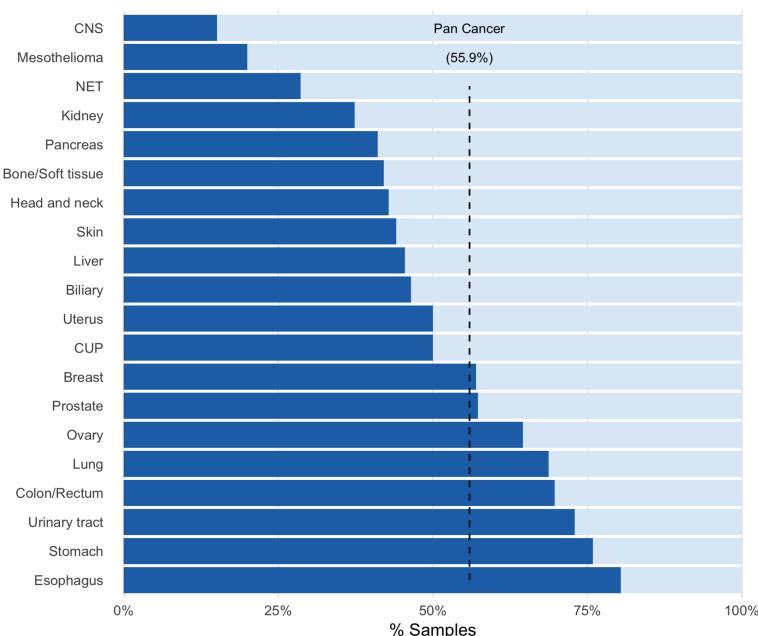
### B. CNS



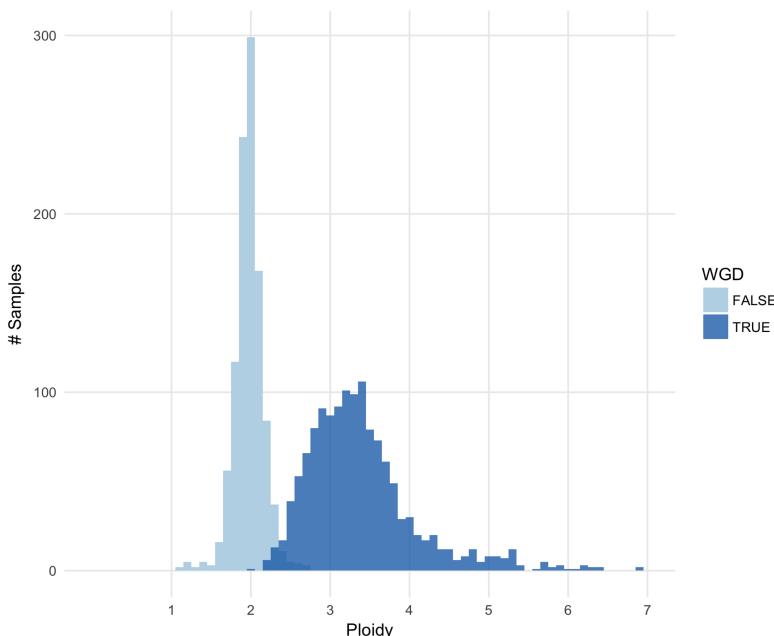
### C. Kidney



### D.

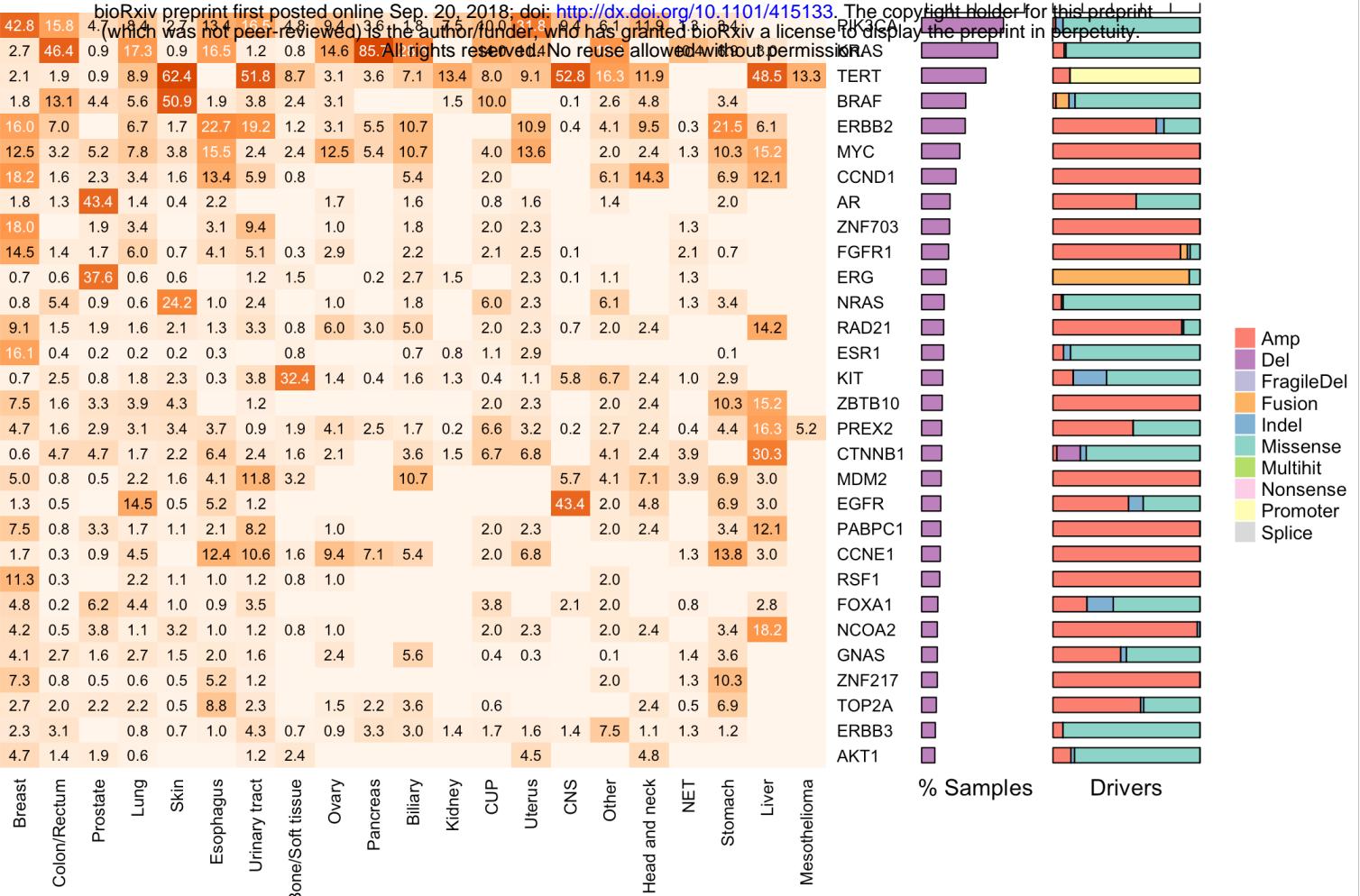


### E.

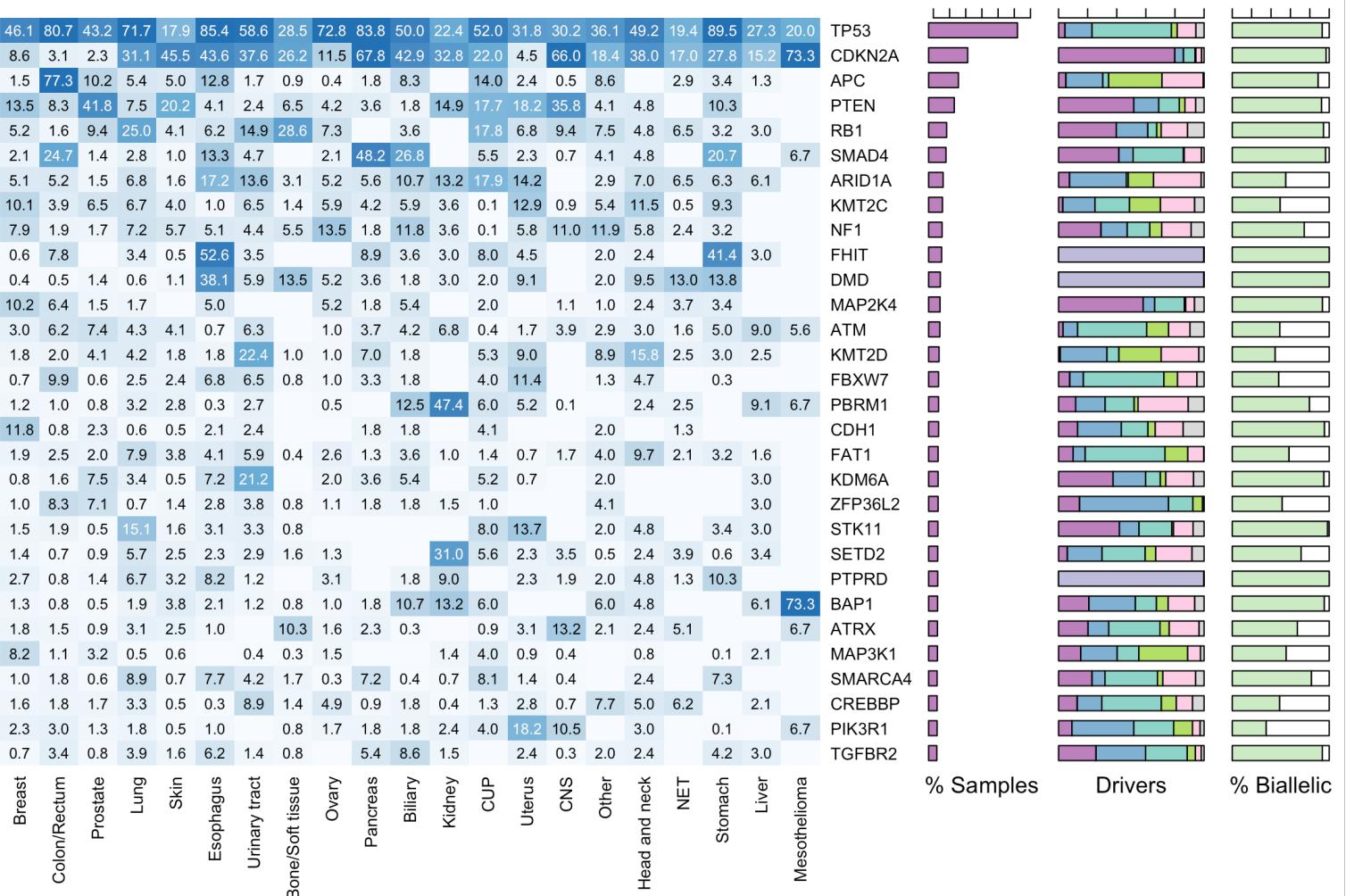


A

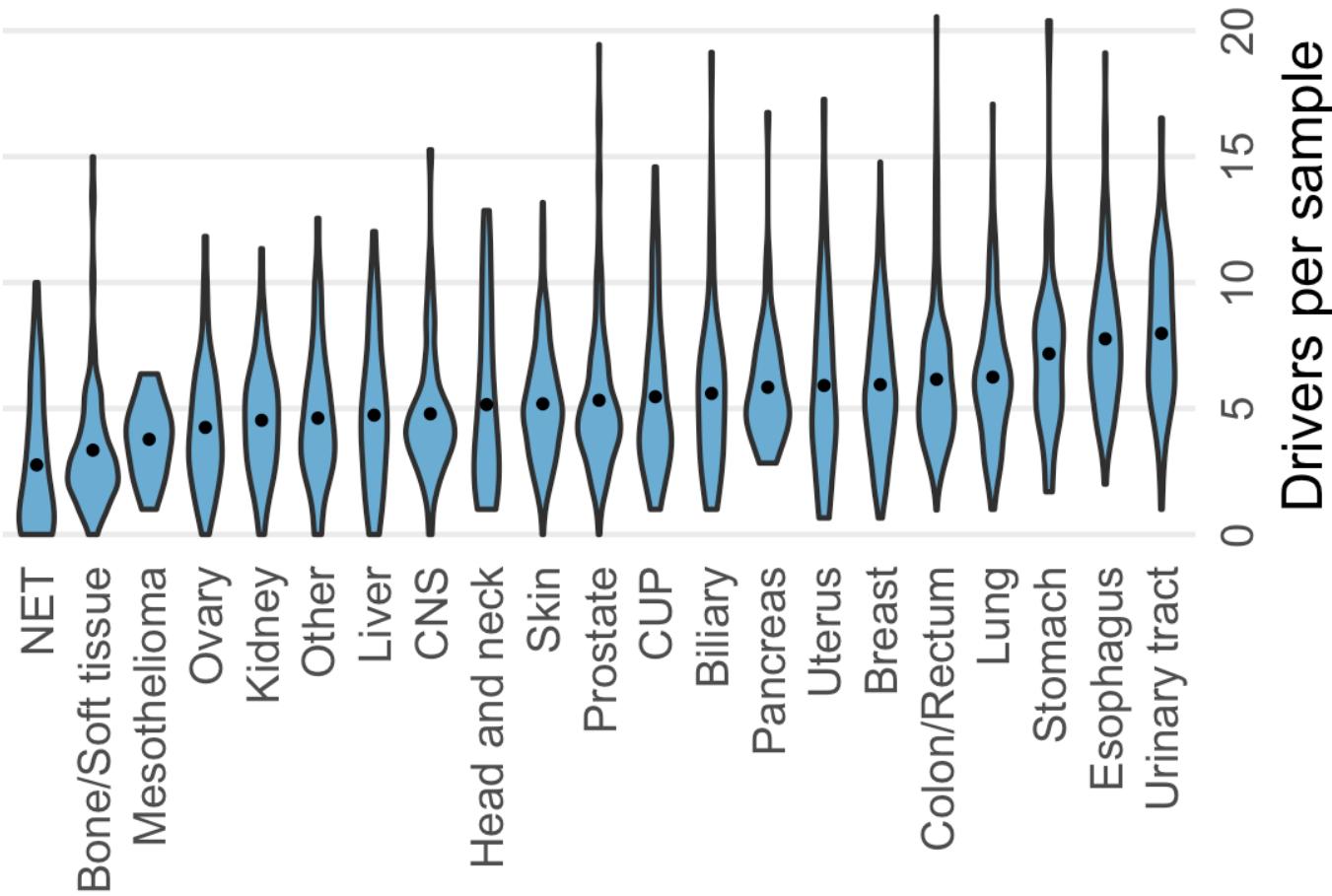
bioRxiv preprint first posted online Sep. 20, 2018; doi: <http://dx.doi.org/10.1101/415133>; this version posted September 20, 2018. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity.



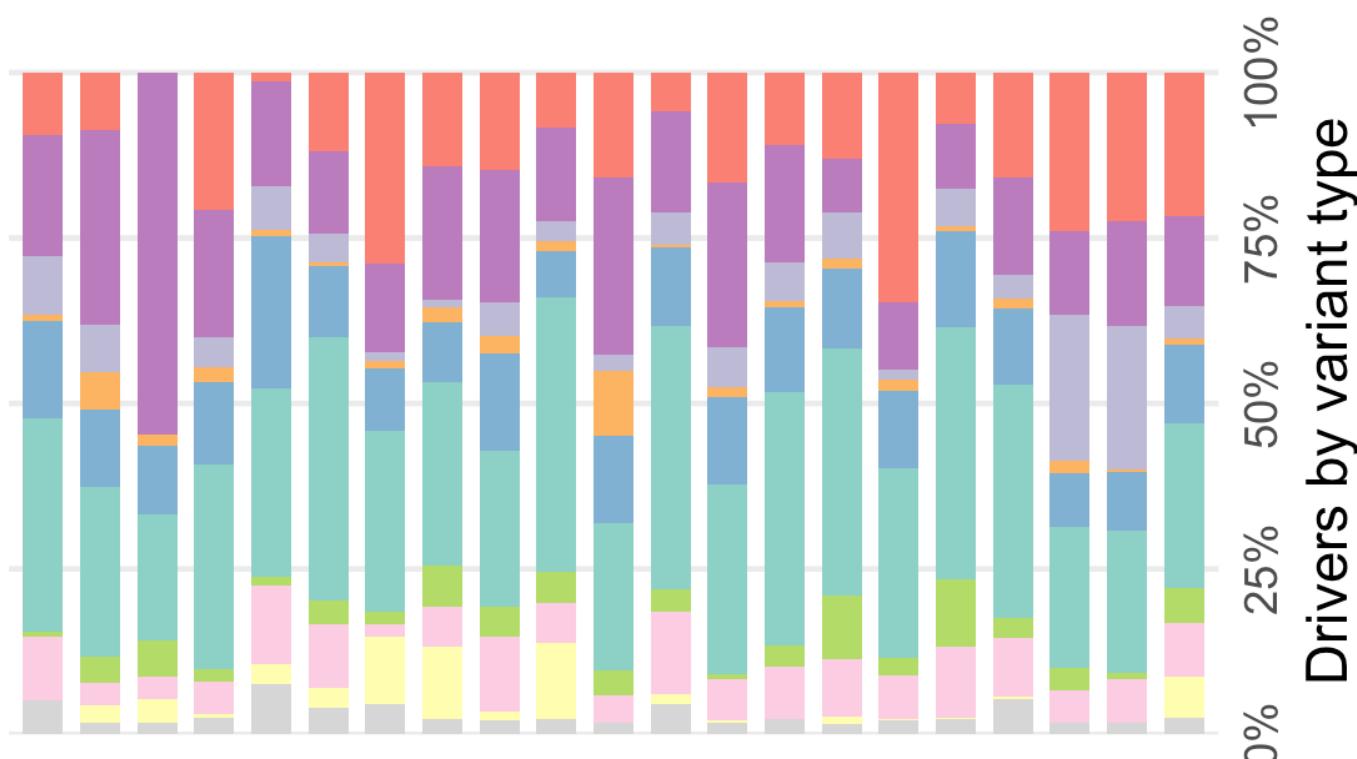
B



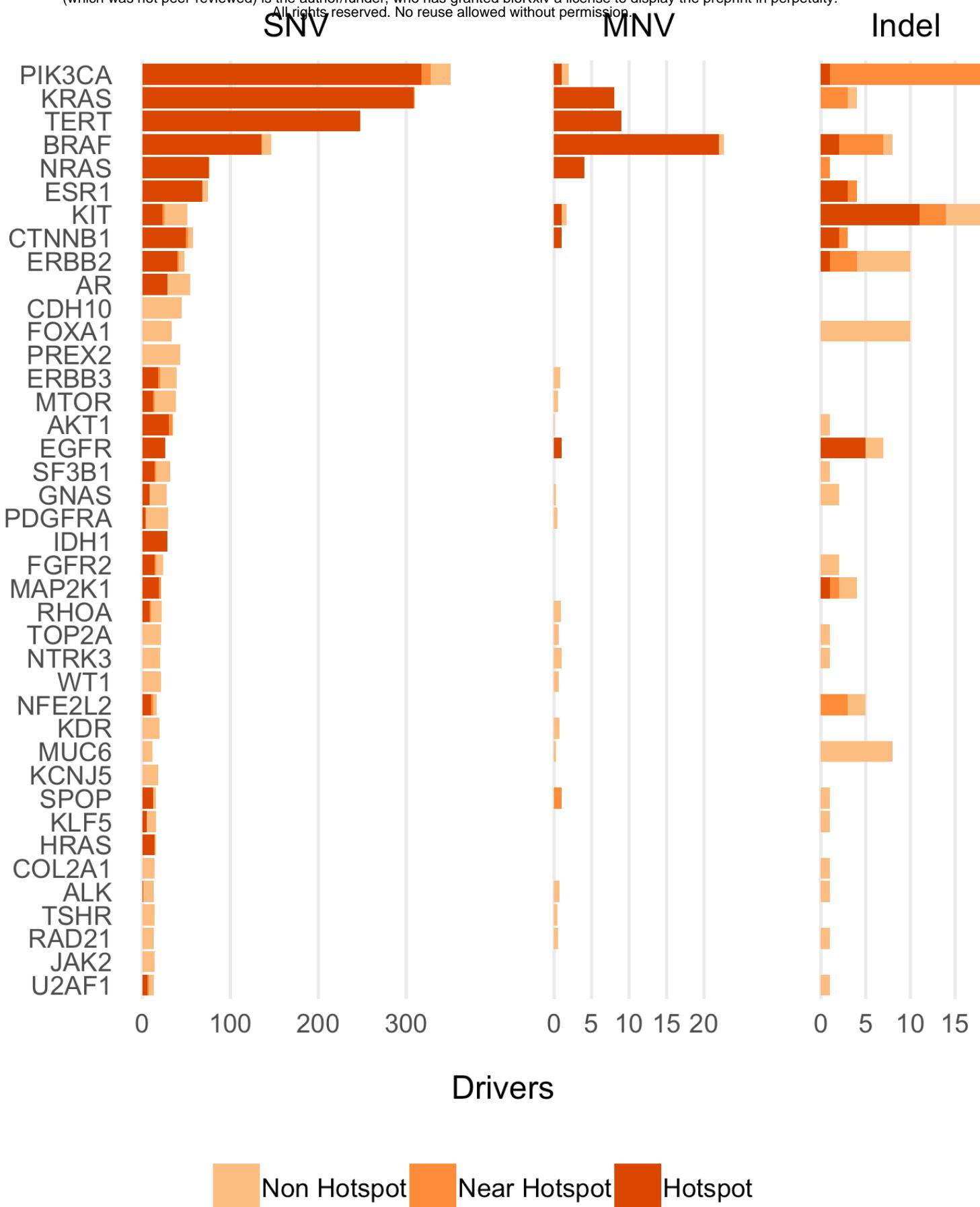
**A**

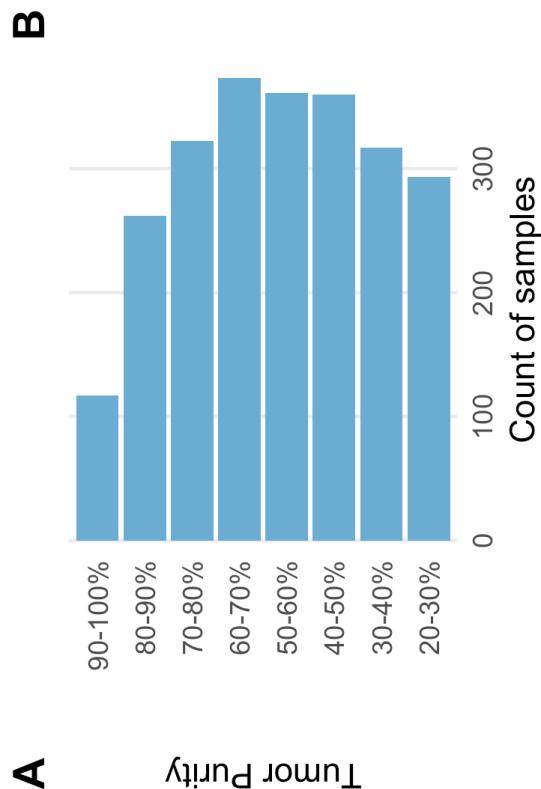
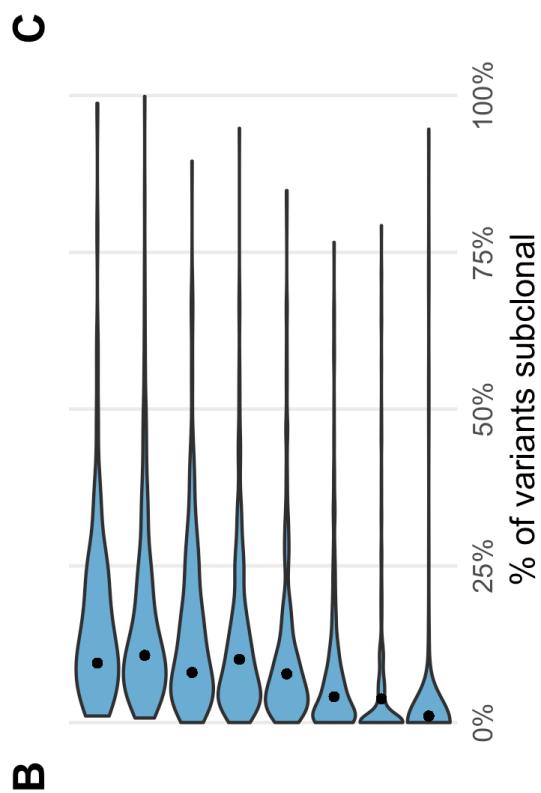
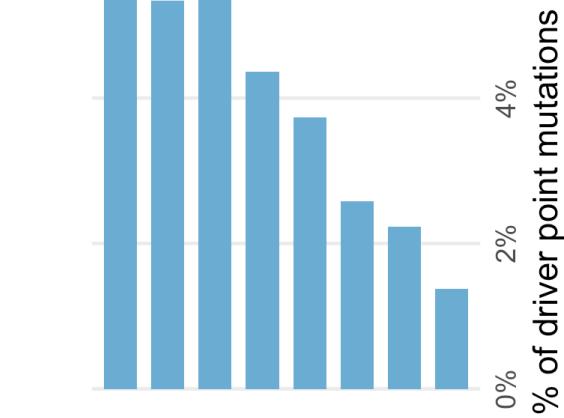


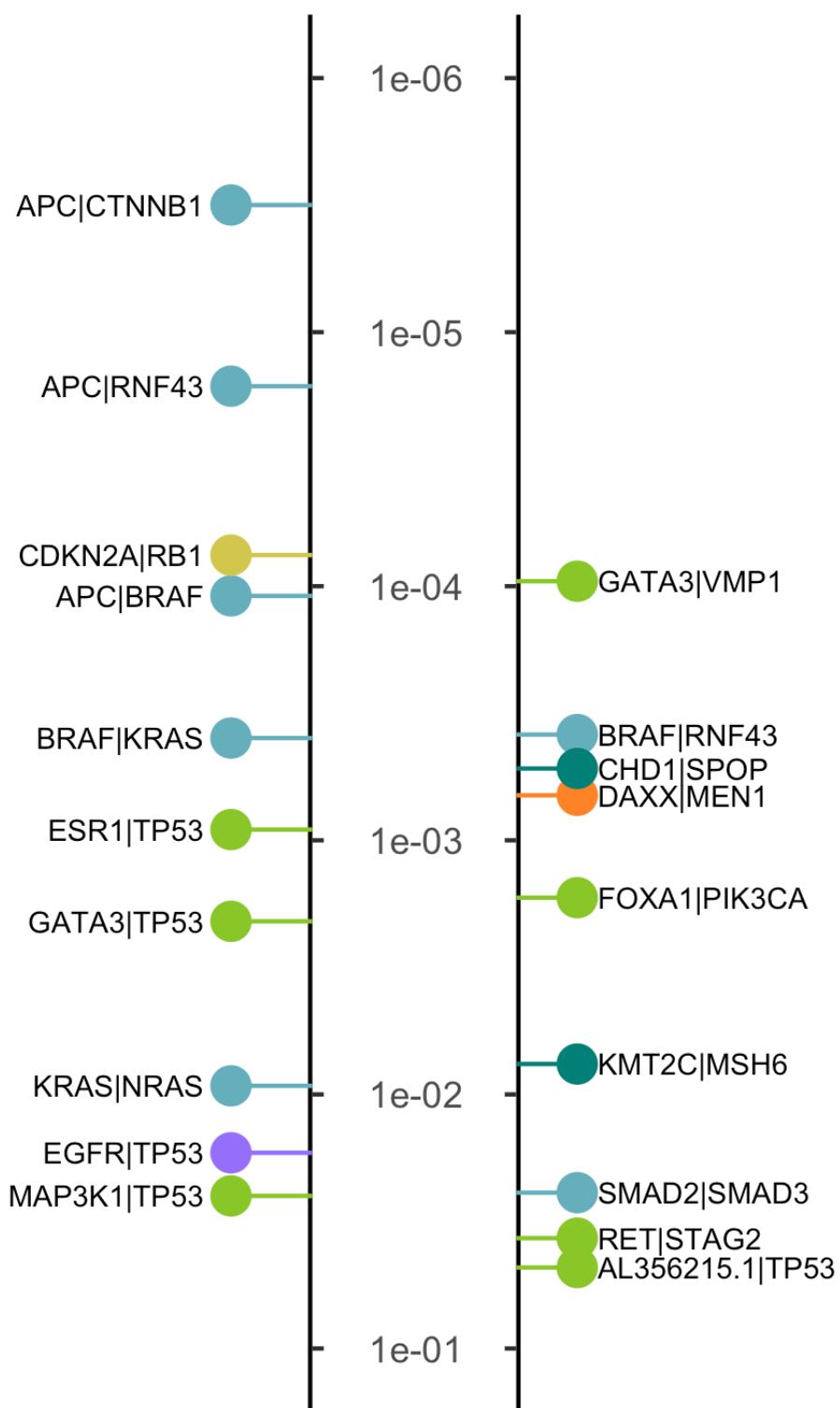
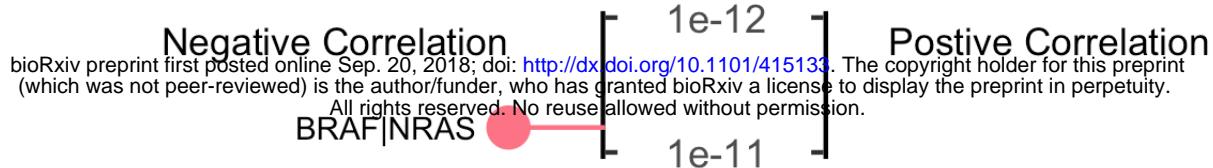
**B**



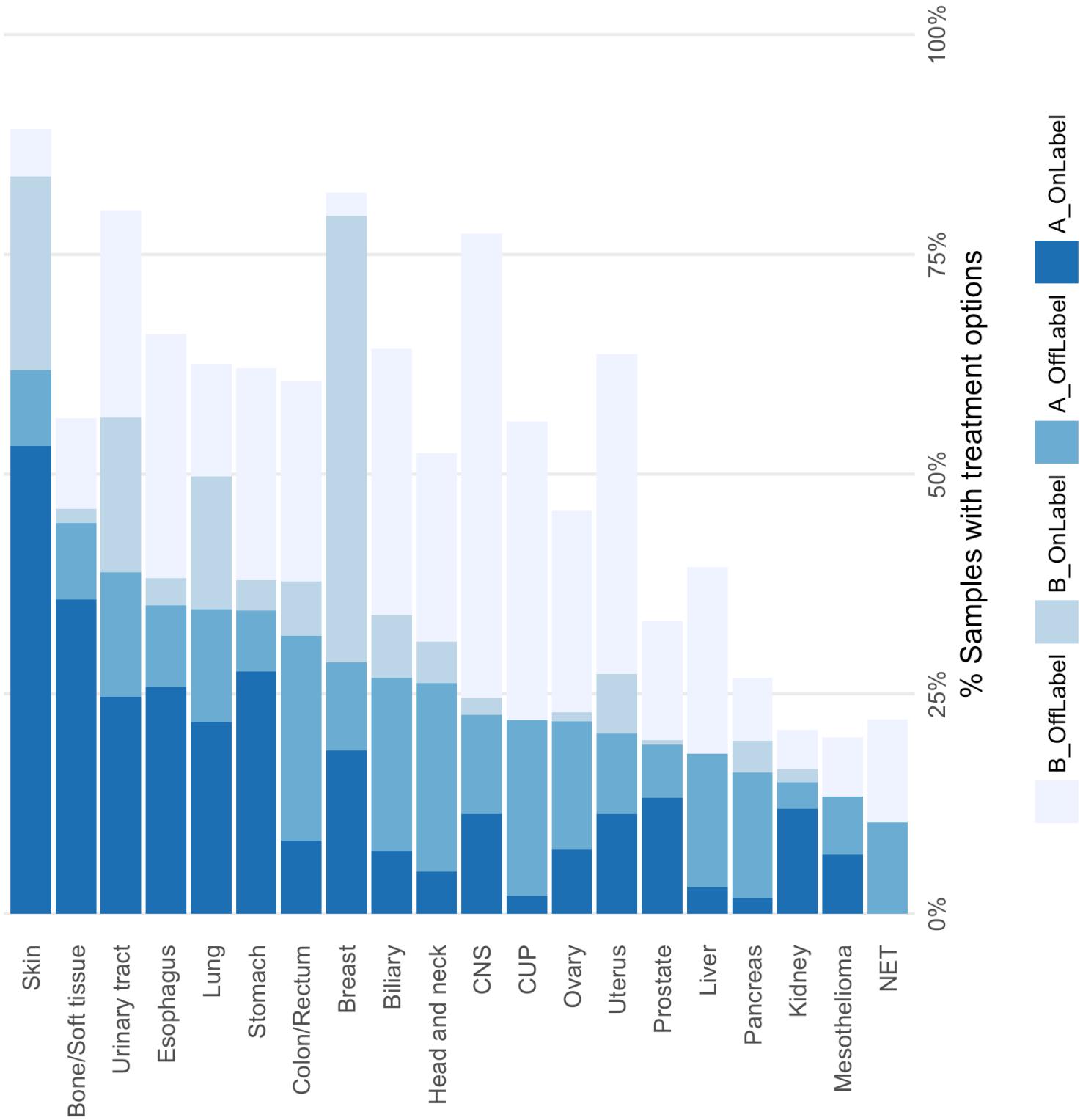
A







- Breast   ● Colon/Rectum   ● NET   ● Skin
- CNS   ● Lung   ● Prostate



# Pan-cancer whole genome analyses of metastatic solid tumors

Peter Priestley, Jonathan Baber, Martijn P. Lolkema, Neeltje Steeghs, Ewart de Bruijn, Korneel Duyvesteyn, Susan Haidari, Arne van Hoeck, Wendy Onstenk, Paul Roepman, Charles Shale, Mircea Voda, Haiko J. Bloemendaal, Vivianne C.G. Tjan-Heijnen, Carla M.L. van Herpen, Mariette Labots, Petronella O. Witteveen, Egbert F. Smit, Stefan Sleijfer, Emile E. Voest, Edwin Cuppen

## Detailed methods

1. Sample collection	2
2. Sequencing workflow	2
3. Somatic point mutation calling	2
4. MNV correction	4
5. Somatic structural variant calling	5
6. Identification of gene fusions	5
7. Copy number calling	6
8. Sample filtering based on copy number output	9
9. Assessment of impact of sequencing depth coverage on variant calling sensitivity	9
10. Clonality and biallelic status of point mutations	9
11. WGD status determination	11
12. MSI status determination	12
13. Holistic gene panel for driver discovery	13
14. Significantly mutated driver genes discovery	13
15. Significantly amplified & deleted driver gene discovery	14
16. Fragile site annotation	15
17. Driver catalog construction	16
18. Driver co-occurrence analysis	19
19. Actionability analysis	19
20. Data availability	22
21. References	23

## 1. Sample collection

Patients with advanced cancer not curable by local treatment options and being candidates for any type of systemic treatment and any line of treatment were included as part of the CPCT-02 (NCT01855477) and DRUP (NCT02925234) clinical studies, which were approved by the medical ethical committees (METC) of the University Medical Center Utrecht and the Netherlands Cancer Institute, respectively. A total of 41 academic, teaching and general hospitals across the Netherlands participated in these studies and collected material and clinical data by standardized protocols<sup>1</sup>. Patients have given explicit consent for whole genome sequencing and data sharing for cancer research purposes. Clinical data, including primary tumor type, biopsy location, gender and birth year were collected in electronic case record forms and stored in a central database.

Core needle biopsies were sampled from the metastatic lesion, or when considered not feasible or not safe, from the primary tumor site when still in situ. One to four biopsies were collected (average of 2.1 per patient) and frozen in liquid nitrogen directly after sampling and further processed at a central pathology tissue facility. Frozen biopsies were mounted on a microtome in water droplets for optimal preservation of all types of biomolecules (DNA, RNA and proteins) for subsequent and future omics-based analyses. A single 6 micron section was collected for hematoxylin-eosin (HE) staining and estimation of tumor cellularity by an experienced pathologist. Subsequently, 25 sections of 20 micron, containing an estimated 25,000 to 500,000 cells, were collected in a tube for DNA isolation. In parallel, a tube of blood was collected in CellSave (Menarini-Silicon Biosystems) tubes, which was shipped by room temperature to the central sequencing facility at the Hartwig Medical Foundation. Left-over material (biopsy, DNA) after sample processing was stored in biobanks associated with the studies at the University Medical Center Utrecht and the Netherlands Cancer Institute.

## 2. Sequencing workflow

DNA was isolated from biopsy and blood on an automated setup (QiaSymphony) according to supplier's protocols (Qiagen) using the DSP DNA Midi kit for blood and QIAsymphony DSP DNA Mini kit for tissue and quantified (Qubit). Typically, DNA yield for the tissue biopsy ranged between 50 and 5,000 ng. A total of 50 - 200 ng of DNA was used as input for TruSeq Nano LT library preparation (Illumina), which was performed on an automated liquid handling platform (Beckman Coulter). DNA was sheared using sonication (Covaris) to average fragment lengths of 450 nt. Barcoded libraries were sequenced as pools (blood control 1 lane equivalent, tumor 3 lane equivalents) on HiSeqX (V2.5 reagents) generating 2 x 150 read pairs using standard setting (Illumina).

BCL output from the HiSeqX platform was converted using Illumina bcl2fastq tool (versions 2.17 to 2.20 have been used) using default parameters. Reads were mapped to the reference genome GRCH37 using BWA-mem v0.7.5a<sup>2</sup>, duplicates were marked for filtering and INDELS were realigned using GATK v3.4.46 IndelRealigner. GATK Haplotype Caller v3.4.46 was run to call germline variants in the reference sample. For somatic SNV and INDEL variant calling, GATK BQSR<sup>3</sup> is also applied to recalibrate base qualities.

## 3. Somatic point mutation calling

We called SNV & INDEL somatic variants using Strelka v1.0.14<sup>4</sup> with the following optimisations:

- **Preservation of known variants:** From the raw Strelka output we marked all known pathogenic variants from external databases such that these would be preserved from all subsequent filtering. The list of pathogenic variants used was the union of:

- Point mutations in CIViC<sup>5</sup> with level C evidence or higher (download = 01-mar-2018)
- Somatic variants from CGI<sup>6</sup> (update: 17-jan-2018)
- Oncogenic or likelyOncogenic variants from OncoKb<sup>7</sup> (download = 01-mar-2018); <http://oncokb.org/api/v1/utils/allAnnotatedVariants.txt>
- TERT promoter variants at genomic coordinates: 5:1295242, 5:1295228, 5:1295250
- **Modified quality score filtering**
  - We split variants into high confidence (HC) and low confidence (LC) regions using the NA12878 GIABv3.2.2 high confidence region definitions<sup>8</sup>, based on the observation that we produce far higher rates of false positives variant calls in LC regions
  - Set quality score cutoffs for SNV & INDEL to 10 for HC regions and 20 for LC regions (default = 15 for SNV, 30 for INDEL)
  - Added an additional quality filter to tighten filtering for low allelic frequency variants: quality score \* allele frequency > 1.3
- **Improved repeat sensitivity:** Switched off the default Strelka repeat filter to improve indel calling in microsatellites and short repeats.
- **Panel of normals (PON) to remove germline leakage:** Filtered out any variants which were found by GATK haplotypecaller in more than 5 samples in a germline PON consisting of 2000 of our reference blood samples. PON available at (<https://resources.hartwigmedicalfoundation.nl/>)
- **PON to remove strelka-specific artefacts:** Filtered any variant which was supported by 2 or more reads in strelka in the reference sample in at least 4 patients in our cohort. PON available at (<https://resources.hartwigmedicalfoundation.nl/>)
- **Removal of INDELS near a PON filtered INDEL** - Regions of complex haplotype alterations are often called as multiple long indels which can make it more difficult to construct an effective PON, and sometimes we find residual artefacts at these locations. Hence we also filter inserts or deletes which are 3 bases or longer where there is a PON filtered INDEL of 3 bases or longer within 10 bases in the same sample.

The settings and tools for this optimized HMF pipeline are available at <https://github.com/hartwigmedical/>. We tested the default and HMF optimized settings on a GIAB mix-in sample (ref = NA24385, tumor = 70% NA24385, 30% NA12878) to test sensitivity at a realistic purity and on a null tumor (ref = NA12878, tumor = NA12878) to test precision. The results of this analysis are as follows:

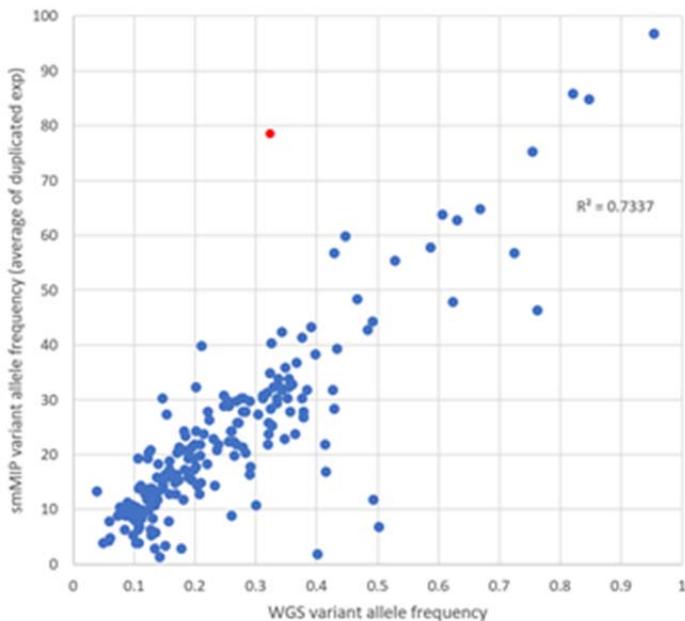
Configuration	SNV sensitivity	SNV false positive / genome	INDEL sensitivity	Indel false positive / genome
Default	93%	3500	24%	41
Optimized HMF pipeline	96%	109	77%	27

We performed external validation of a set of single nucleotide variants (SNV) and short insertion/deletions (indels) that have been detected by Whole Genome Sequencing (WGS) using the single molecule Molecular Inversion Probe (smMIP) technology<sup>9</sup>. SNV and short indels variants were semi-randomly selected from 30 patient samples. The first selection was to include every variant that was reported in the patient report (114 gene panel) as obtained during the routine CPCT-02 study analysis. This way, a total of 82 variants (67 SNVs, 15 indels) were selected in 45 genes. The second selection involved random sampling of 256 coding and non-coding variants from the same 30 patient samples.

A custom smMIP panel was designed to cover the selected variants. For 45 variants (17.6%) no smMIP design was possible, all of which were intergenic variants. For the other 211 variants probes could successfully be designed. Analysis of the smMIP sequencing data indicated that for 17 of the 211 variants

(8.1%) the smMIP sequencing data was of insufficient quality (mostly due to repeat stretches), while the WGS data seemed sufficiently reliable for accurate calling (confirmed by visual inspection of the read data), including 3 coding variants (*RB1*, *ERBB4* and *BRCA2*) and 14 intergenic regions. The retrospective investigation of the WGS data indicated that for another three variants (1.4%) the smMIP as well as the WGS data was of insufficient quality due to large homopolymer stretches.

In total 192 variants could be successfully sequenced and analyzed using the smMIP and could be used for confirmation of the WGS findings. 189 SNVs and indel variants (98.4%) were confirmed by smMIP sequencing, indicating a very high accuracy of WGS-derived variant calling results. All three variants that could not be confirmed by smMIP were from intergenic regions, including 1 variant that showed a mixed double-variant (chr3:75887550\_G>T/C) and for which both technologies had difficulties in accurately calling the genotype. For the remaining 2 variants (chr8:106533360\_106533361insAC, chr12:125662751\_125662752insA), it remains unclear if these could not be detected by smMIP or were falsely called by WGS, as they fall in repetitive genomic stretches.



The 189 successfully confirmed variants showed a good linear correlation in variant allele frequency between WGS and smMIP sequencing (average of duplicates) with an  $R^2$  of 0.733. This result indicated that WGS, with its lower read depth (on average between 100-110x) than smMIP and without a read-barcoding system, is accurate in quantitatively determining the variant frequency at frequencies above 5%. One variant (ch19:55276095C>T, indicated in red in the figure above) showed a large deviation in variant frequency, which was likely caused due to the much lower than expected coverage of the variant, both in the WGS (37 reads) as well as in the smMIP data (28 and 35 reads).

#### 4. MNV correction

Strelka somatic variants that appear on consecutive positions, or 1 base apart were considered potential multi nucleotide variants (MNVs). The BAM files were re-examined, and the variants were merged into a single MNV if greater than 80% of the reads with a mapping quality score of at least 10 and which are neither unmapped, duplicated, secondary, nor supplementary containing any of the individual variants also contained the other variants of the potential MNV. The attributes of the resulting MNV

variant were determined by picking the minimum values from the individual variants forming the MNV. MNVs were marked as PON filtered only if both individual variants were PON filtered.

## 5. Somatic structural variant calling

Structural Variants were called using Manta(v1.0.3)<sup>10</sup> with default parameters. We then re-examined each breakpoint calculated variant allele frequencies for each break end and applied seven additional filters to the Manta output to improve precision using an internally built tool called 'Breakpoint-Inspector' (BPI) v1.5. Two main types of filters are applied by BPI:

- **Evidence of variant in reference sample** - variants are filtered if we can find any evidence of paired read support , split read support or soft clipping concordance (5+ bases at exact breakpoint) in the matching blood sample.
- **Inadequate support for variant in tumor sample** - For all inversions and translocations and for long deletions and tandem duplications (>1000 bases between breakpoints) we require at least 1 read with paired read support. For short deletions and duplications (<1000 bases between breakpoints) we require at least 1 read with split read support. In both cases at least one of those reads must be anchored with at least 30 bases at each breakpoint. We also require the minimum read coverage across each breakpoint in the tumor to be > 10 depth.

Code and description of filters for BPI are available at

<https://github.com/hartwigmedical/hmftools/tree/master/break-point-inspector>.

Each break end was annotated with it's position in all transcripts from 'KNOWN' genes in Ensembl v89.37<sup>11</sup>. Each gene was marked as disrupted if there was at least one structural variant that impacted on the canonical transcript.

## 6. Identification of gene fusions

For each structural variant, every combination of annotated overlapping transcripts from each breakend was tested to see if it could potentially form an intronic inframe fusion. A list of 411 curated known fusion pairs was sourced by taking the union of known fusions from the following external databases:

- Cosmic curated fusions<sup>12</sup> (v83)
- OncoKb<sup>7</sup> (download = 01-mar-2018)
- CGI<sup>6</sup> (update: 17-jan-2018)
- CIViC<sup>5</sup> (download = 01-mar-2018)

We then also created a list of promiscuous fusion partners using the following rules

- **3' promiscuous:** Any gene which appears on the 3' side in more than 3 of the curated fusion pairs OR appears at least once on the 3' side and is marked as promiscuous in either OncoKb, CGI or CIViC
- **5' promiscuous:** Any gene which appears on the 5' side in more than 3 of the curated fusion pairs OR appears at least once on the 5' side and is marked as promiscuous in either OncoKb, CGI or CIViC

For each promiscuous partner we also curated a list of essential domains that must be preserved to form a viable fusion partner.

Finally, we report an intronic inframe fusion if the following conditions are met

- Matches an exact fusion from the curated list OR is intergenic and matches 5' promiscuous OR matches 3' promiscuous gene

- Curated domains are preserved
- Does not involve the 3'UTR region of either gene
- For intragenic fusions, must start and end in coding regions of the gene

## 7. Copy number calling

We use an in house developed integrated tool, PURity & PLoidy Estimator (PURPLE), that combines B-allele frequency (BAF), read depth and structural variants to estimate the purity and copy number profile of a tumor sample. Version v2.14 of PURPLE has been used.

There are 5 key steps in the PURPLE pipeline:

### 1. Calculate BAF in tumor at high confidence heterozygous germline loci

We determine the BAF of the tumor sample by finding heterozygous locations in the reference sample from a panel of 796,447 common germline heterozygous SNP locations. To ensure that we only capture heterozygous points, we filter the panel to only loci with allelic frequencies in the reference sample between 40% and 60% and with depth between 50% and 150% of the reference sample genome wide average. Typically, this yields 140k-200k heterozygous germline variants per patient. We then calculate the allelic frequency of corresponding locations in the tumor.

### 2. Determine read depth ratios for tumor and reference genomes

The raw read counts per 1,000 base window for both normal and tumor samples, by counting the number of alignment starts in the respective bam files with a mapping quality score of at least 10 that is neither unmapped, duplicated, secondary, nor supplementary. Windows with a GC content less than 0.2 or greater than 0.6 or with an average mappability below 0.85 are excluded from further analysis.

Next we apply a GC normalization to calculate the read ratios. We divide the read count of each window by the median read count of all windows sharing the same GC content then normalise further to the ratio of the median to mean read count of all windows.

Finally, the reference sample ratios have a further ‘diploid’ normalization applied to them to remove megabase scale GC biases. This normalization assumes that the median ratio of each 10Mb window (minimum 1Mb readable) should be diploid for autosomes and haploid for sex chromosomes in males in the germline sample.

### 3. Segmentation

We segment the genome into regions of uniform copy number by combining segments generated from the read ratios for both tumor and reference sample, from the BAF points with structural variant breakpoints derived from Manta & BPI. Read ratios and BAF points are segmented independently using the Bioconductor copynumber package<sup>13</sup> which uses a piecewise constant fit (PCF) algorithm (with custom settings gamma = 100, k = 1). These segment breaks are then combined with the structural variants breaks according to the following rules:

1. Every structural variant break starts a new segment, as does chromosome starts, ends and centromeres. This is regardless of if they are distinguishable from existing segments or not.
2. Ratio and BAF segment breaks are only included if they are distinguishable from an existing segment.
3. To be distinguishable, a break must be at least one complete mappable read depth window away from an existing segment.

Once the segments have been established we map our observations to them. In each segment we take the median BAF of the tumor sample and the median read ratio of the tumor and reference samples. We also record the number of BAF points within the segment as the BAFCOUNT.

A reference sample copy number status is determined at this stage based on the observed copy number ratio in the reference sample, either 'DIPLOID' ( $0.8 \leq \text{read depth ratio} \leq 1.2$ ), 'HETEROZYGOUS\_DELETION' ( $0.1 \leq \text{ratio} < 0.8$ ), 'HOMOZYGOUS\_DELETION' ( $\text{ratio} < 0.1$ ), 'AMPLIFICATION' ( $1.2 < \text{ratio} \leq 2.2$ ) or 'NOISE' ( $\text{ratio} > 2.2$ ). The purity fitting and smoothing steps below use only the DIPLOID germline segments.

#### 4. Purity Fitting

Next we jointly fit tumor purity and sample ploidy (expressed as a normalisation factor) according to the following principles:

1. The absolute copy number of each segment should be close to an integer ploidy
2. The BAF of each segment should be close to a % implied by integer major and minor allele ploidies.
3. Higher ploidies have more degenerate fits but are less biologically plausible and should be penalised
4. Segments are weighted by the count of BAF observations which is treated as a proxy for confidence of BAF and read depth ratio inputs.
5. Segments with lower observed BAFs have more degenerate fits and are weighted less in the fit

For any given tumor purity and sample ploidy we calculate the score by first modelling the major and minor allele ploidy of each segment and minimising the deviation between the observed and modelled values according to the following formulas:

ModelDeviation = abs(ObservedRatio - ModelRatio) + abs(ObservedBaf - ModelBaf)  
ModelBaf = (tumorPurity \* (segmentMinorPloidy - 1) + 1) / (tumorPurity \* (segmentPloidy - 2) + 2)  
ModelRatio = sampleNormFactor + (segmentPloidy - 2) \* tumorPurity \* sampleNormFactor / 2d;

Once modelled, each segment is given a ploidy penalty:

PloidyPenalty = 1 + min(SingleEventDistance, WholeGenomeDoublingDistance);  
WholeGenomeDoublingDistance = 1 + abs(segmentMajorAllele - 2) + abs(segmentMinorAllele - 2);  
SingleEventDistance = abs(segmentMajorAllele - 1) + abs(segmentMinorAllele - 1);

Summing up over all the segments generates a score for each tumor purity / sample ploidy combination from which we can select the minimum:

$$\begin{aligned} \text{FittedPurityScore} \\ = \frac{1}{\text{TotalBafCount}} \sum_{i=1}^n & \text{PloidyPenalty}_i \times \text{ModelDeviation}_i \times \text{BafCount}_i \\ & \times \text{ObservedBaf}_i \end{aligned}$$

Given a fitted purity and sample ploidy we are then able to determine the purity adjusted copy number and BAF of each segment in the tumor from the unadjusted read ratios and BAFs respectively.

The purity estimates of PURPLE were validated using the tumor cell line COLO829. We created diluted in-silico mixture models of the tumor and blood cell lines for COLO829 with simulated purities of 20%, 30%, 40%, 60%, 80% and 100%, and ran PURPLE on the simulated BAM files against a reference sample.

The PURPLE estimates match the simulation closely as follows:

Simulated Purity	PURPLE estimated purity	Difference
20%	20%	0%
30%	30%	0%
40%	40%	0%
50%	50%	0%
60%	60%	0%
80%	81%	1%
100%	100%	0%

## 5. Smoothing

Since the segmentation algorithm is highly sensitive, and there is a significant amount of noise in the read depth in whole genome sequencing, many adjacent segments created above will have a similar copy number and BAF profile and can be combined and averaged to form a larger, smoothed, region.

We apply a number of rules to merge adjacent regions to create a smooth copy number profile.

1. Never merge a segment break created from a structural variant break end.
2. Use the count of BAF points as a proxy for confidence or weight in the region. Note that some segments may have a BAF count of 0.
3. Merge segments where the difference in BAF and copy number is within tolerances.
4. BAF tolerance is linear between 0.03 and 0.35 dependent on BAF count.
5. Copy number tolerance is linear between 0.3 and 0.7 dependent on BAF count. The tolerance also increases linearly as purity of the tumor sample decreases below 20%.
6. Start from most confident segment and smooth outwards until we reach a segment outside of tolerance. Move on to next highest unsmoothed section.
7. It is possible to merge in (multiple) segments that would otherwise be outside of tolerances if:
  - a. The total dubious region is sufficiently small (<30k bases or <50k bases if approaching centromere); and
  - b. The dubious region does not end because of a structural variant; and
  - c. The dubious region ends at a centromere, telomere or a segment that is within tolerances.
8. When the entire short arm of a chromosome is lacking copy number information (generally on chromosome 13, 14, 15, 21, or 22), the copy number of the long arm is extended to the short arm.
9. Any remaining unknown segments are given the expected copy number of their associated chromosome, i.e. 2 for autosomes and female allosomes, 1 for male allosomes.

Where clusters of SVs exist which are closer together than our read depth ratio window resolution of 1,000 bases, the segments in between will not have any copy number information associated with them.

To resolve this, we infer the ploidy from the surrounding copy number regions. The outermost segment of any SV cluster will be associated with a structural variant with a ploidy that can be determined from the adjacent copy number region and the VAF of the SV. We use this ploidy and the orientation of structural variant to calculate the change in copy number across the SV and hence the copy number of the outermost unknown segment. We repeat this process iteratively and infer the copy number of all regions within a cluster.

Once region smoothing is complete, it is possible there will be regions of unknown BAF, if no BAF points were present in a copy number region. We infer this BAF by assuming that they share their minor allele ploidy with their neighbouring region. If there are multiple neighbouring regions with known BAF we use the highest confident region (i.e. highest BAF count) to infer.

At this stage we have determined a copy number and minor allele ploidy for every base in the genome.

## **8. Sample filtering based on copy number output**

Following our copy number calling, samples were QC filtered from the analysis based on 4 criteria:

- **NO\_TUMOR** - If PURPLE fails to find any aneuploidy AND the number of somatic SNVs found is less than 1,000 then the sample is marked as NO\_TUMOR.
- **MIN\_PURITY** - We exclude samples with a purity of <20%
- **FAIL\_SEGMENT** - We remove samples with more than 120 copy number segments unsupported at either end by SV breakpoints. This step was added to remove samples with extreme GC bias, with differences in depth of up to or in excess of 10x between high and low GC regions. GC normalisation is unreliable when the corrections are so extreme so we filter.
- **FAIL\_DELETED\_GENES** - We removed any samples with more than 280 deleted genes. This QC step was added after observing that in a handful of samples with high MB scale positive GC bias we sometimes systematically underestimate the copy number in high GC regions. This can lead us to incorrectly infer homozygous loss of entire chromosomes, particularly on chromosome 17 and 19.

Where multiple biopsies exist for a single patient, we always choose the highest purity sample for our analysis of mutational load and drivers.

## **9. Assessment of impact of sequencing depth coverage on variant calling sensitivity**

To assess the impact of our sequencing depth on variant calling sensitivity, we selected 10 samples at random, downsampled the BAMs by 50%. We then reran the identical somatic variant calling pipeline.

Comparing the output to the original runs, we found near identical purities and ploidies for the down sampled runs ([Extended Data Fig. 2](#)). We observed an average decrease in sensitivity of 10% for SNV, 15% for MNV, 19% for SV, and 2% for INDEL.

## **10. Clonality and biallelic status of point mutations**

For each point mutation we determined the clonality and biallelic status by comparing the estimated ploidy of the variant to the local copy number at the exact base of the variant. The ploidy of each variant is calculated by adjusting the observed VAF by the purity and then multiplying by the local copy number to work out the absolute number of chromatids that contain the variant.

We mark a mutation as biallelic (i.e. no wild type remaining) if Variant Ploidy > Local Copy Number - 0.5. The 0.5 tolerance is used to allow for the binomial distribution of VAF measurements for each variant. For example, if the local copy number is 2 than any somatic variant with measured ploidy > 1.5 is marked as biallelic.

For each variant we also determine a probability that it is subclonal. This is achieved via a 2 step process

### **1. Fit the somatic ploidies for each sample into a set of clonal and subclonal peaks**

We apply an iterative algorithm to find peaks in the ploidy distribution:

- Determine the peak by finding the highest density of variants within +/- 0.1 of every 0.01 ploidy bucket.
- Sample the variants within a 0.05 ploidy range around the peak.
- For each sampled variant, use a binomial distribution to estimate the likelihood that the variant would appear in all other 0.05 ploidy buckets.
- Sum the expected variants from the peak across all ploidy buckets and subtract from the distribution.
- Repeat the process with the next peak

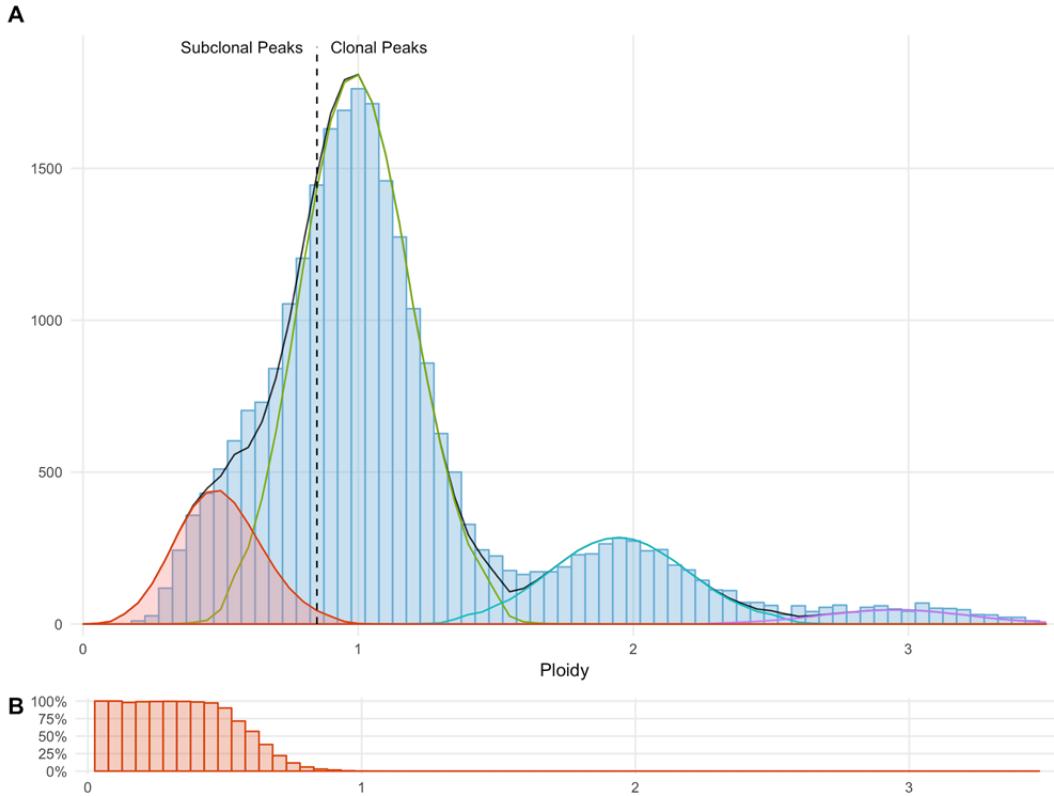
This process yields a set of ploidy peaks, each with a ploidy and a total density (i.e. count of variants). To avoid overfitting small amounts of noise in the distribution, we filter out any peaks that account for less than 40% of the variants in the ploidy bucket at the peak itself. After this filtering we scale the fitted peaks by a constant so that the sum of fitted peaks = the total variant count of the sample.

Finally we mark a peak as subclonal if the peak ploidy < 0.85

### **2. Calculate the probability that each individual variant belongs to each peak**

Once we have fitted the somatic ploidy peaks and determined their clonality, we can calculate the subclonal likelihood for any individual variant as the proportion of subclonal variants at that same ploidy.

The following diagram illustrates this process for a typical sample. Figure A shows the histogram of somatic ploidy for all SNV and INDEL in blue. Superimposed are four peaks in different colours fitted from the sample as described above. The red filled peak is below the 0.85 threshold and is thus considered subclonal. The black line shows the overall fitted ploidy distribution. Figure B shows the likelihood of a variant being subclonal at any given ploidy.



Subclonal counts in this paper are calculated as the total density of the subclonal peaks for each sample. Subclonal driver counts are calculated as the sum across the driver catalog of subclonal probability \* driver likelihood (driver likelihood is explained in detail below).

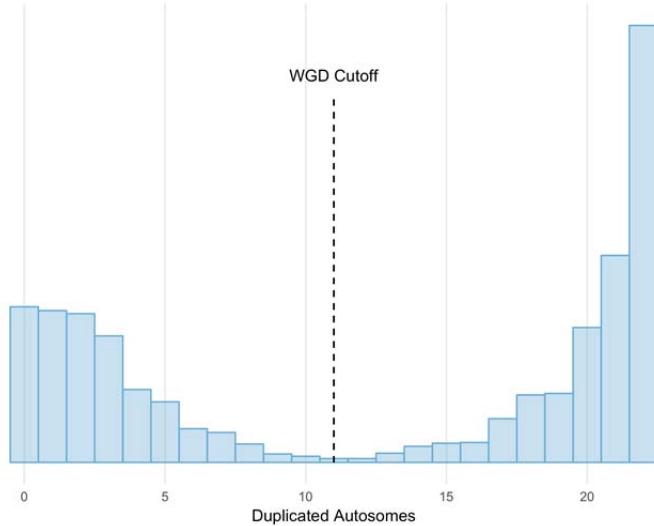
### 11. WGD status determination

We implement a simple heuristic that determines if Whole Genome Duplication has occurred:

#### Major allele Ploidy $>1.5$ on at least 50% of at least 11 autosomes

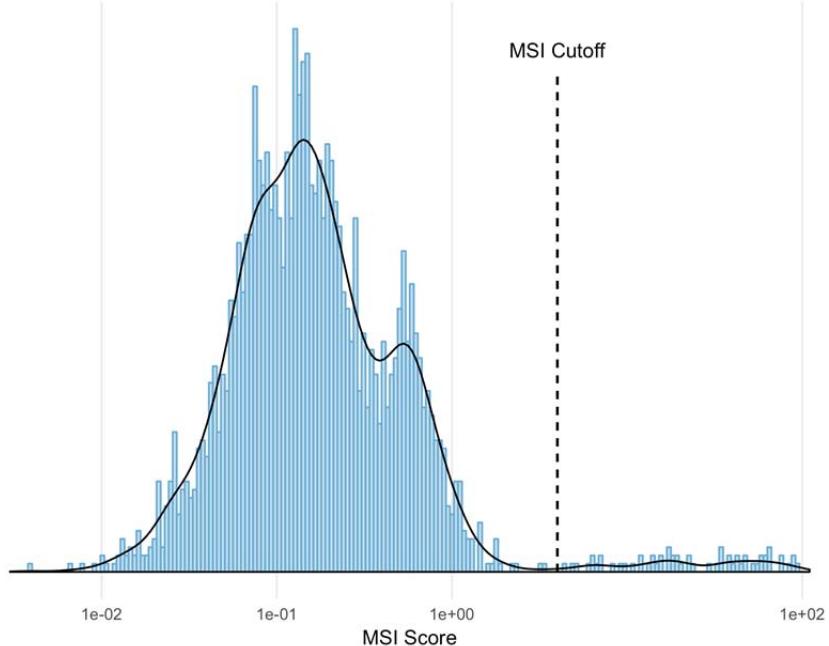
The principle behind this heuristic is that if sufficient independent chromosomes are predominantly duplicated, the most parsimonious explanation is that the duplication occurred in a single genome-wide event.

The number of duplicated autosomes per sample (ie. the number of autosomes which satisfy the above rule) follows a bimodal distribution with 95% of samples have either  $\leq 6$  or  $\geq 15$  autosomes duplicated. Hence, the classification of a genome as WGD is not particularly sensitive to the choice of cut-off as is evident the following chart:



## 12. MSI status determination

To determine the MSI status of all samples we used the method described by the MSIsq tool<sup>14</sup>. In brief, we count the number of INDELS per million bases occurring in homopolymers of 5 or more bases or dinucleotide, trinucleotide and tetranucleotide sequences of repeat count 4 or more. MSIsq scores ranged from 0.004 up to 98.63, with a long tail towards lower MSI scores as shown in the following chart:



To be able to accurately set and validate the MSIsq cutoff for classification of MSI we compared the WGS results with the standard, routinely used MSI assessment using a 5-marker PCR panel (BAT25, BAT26, NR21, NR24 and MONO27 markers). For a batch of 48 pre-selected samples, the MSI PCR assay was blindly performed by an independent ISO-accredited pathology laboratory. Both the binary MSI and MSS classifications were determined, but also the number of positive markers.

A sample was considered as MSI if two or more of the five markers were score as positive (instable). PCR-based analysis identified 16 MSI samples, all of which were also identified by MSIsq with scores >4. MSIsq identified one sample that was missed by PCR-based analysis, although this sample showed microsatellite instability for one out the five markers. The MSIsq scores thus highly correlate with the number of positive MSI PCR markers and all, except one, samples with an elevated score are classified as MSI by pathology. Based on this data we determined the best cutoff for MSIsq classification to be at a **score of 4**.

Results of the PCR-based and WGS based MSI classification are summarized in the table below. The sensitivity of WGS-based MSI classification on this set was 100% (95%CI 82.6 – 100%) with a specificity of 97% (95%CI 88.2-96.9%). The calculated Cohen's kappa score was 0.954 (95%CI 0.696-0.954), indicative of a very high agreement.

	PCR-MSS	PCR-MSI	Total
MSIsq -MSS	31	0	31
MSIsq - MSI	1	16	17
Total	32	16	48

### 13. Holistic gene panel for driver discovery

We used Ensembl release 89 as a basis for our gene definitions and have taken the union of Entrez identifiable genes and protein coding genes as our base panel.

Certain genes have multiple definitions. NPIPA7 for example has two definitions which are equally valid, [ENSG00000214967](#) and [ENSG00000183889](#). To solve this we select a single gene definition based on the following steps:

- 1) Exclude non protein coding genes.
- 2) Favour genes that are present in both Havana and Ensembl.
- 3) Select gene with longest transcript.

This returns our final gene panel tally to 25,963 genes of which 20,083 genes are protein coding. For each gene we chose the canonical transcript or the longest if no canonical transcript exists.

For CDKN2A, we included both the p16 and p14arf transcripts in the analysis given the known importance of both transcripts to tumorigenesis<sup>15</sup> and the fact that the two transcripts use alternate reading frames in the same exon.

### 14. Significantly mutated driver genes discovery

Using all SNV and INDEL variants from the holistic gene panel, we ran dNdScv<sup>16</sup> to find significantly mutated genes (SMGs) and also to estimate the proportion of missense, nonsense, essential splice site and INDEL variants which are drivers in each individual gene in the panel.

Pan cancer and at an individual cancer level we tested the normalised dNdS rates against a null hypothesis that dNdS = 1 for each variant subtype. To identify SMGs in our cohort we used a strict significance cutoff of q<0.01.

2 of the HMF SMG candidates were subsequently removed via manual curation as they were deemed to be likely artefacts of our methods:

- POM121L12 - found only to be significant due to an extreme covariate value in dndscv
- TRIM49B - found to have poor mappability on nearly all its variants and a known close paralog

## 15. Significantly amplified & deleted driver gene discovery

To search for significantly amplified and deleted genes we first calculated the minimum exonic copy number per gene across our holistic gene panel. For amplifications, we searched for all the genes with high level amplifications only (defined as minimum Exonic Copy number  $> 3 * \text{sample ploidy}$ ). For deletions, we searched for all the genes in each sample with either full or partial gene homozygous deletions (defined as minimum exonic copy number  $< 0.5$ ). The Y chromosome was excluded from the deletion analysis since the Y chromosome is deleted altogether in 35% of all male cancer samples in our cohort and hence is difficult to distinguish at the gene level.

We then searched separately for amplifications and deletions, on a per chromosome basis, for the most significant focal peaks, using an iterative GISTIC-like peel off method<sup>17</sup>, specifically:

- Find the highest scoring gene.
  - For deletions the score is simply the count of samples with homozygous deletions in the gene.
  - For amplifications, we need to consider both the count and strength of the amplification so we use:
    - $\text{score} = \text{sum}(\log_2(\text{copy number} / \text{sample ploidy}))$ .
- Record gene as a peak, and mark all consecutive genes with a score within 15% and 25% of the highest score for deletions and amplifications respectively as part of the candidate peak.
- ‘Peel’ off all samples which contributed to the peak across the entire chromosome
- Repeat the process

A filter was applied where we removed deletions from a handful of noisy copy number regions in the genome where we found more than 50% of the observed deletions were not supported on either breakend by a structural variant.

Most of the deletion peaks resolve clearly to a single target gene reflecting the fact that homozygous deletions are highly focal, but for amplifications this is not the case and the majority of our peaks have 10 or more candidates. We therefore annotated the peaks, to choose a single putative target gene using an objective set of automated curation rules in order of precedence:

- If more than 50% of the copy number events in the peeled samples include the telomere or centromere then mark as <CHR>\_<ARM>\_<TELOMERE/CENTROMERE>
- Else choose highest scoring candidate gene which matches a list of actionable amplifications from OncoKB, CGI and CIViC clinical annotation DBs.
- Else choose highest scoring candidate gene found in our panel of significantly mutated genes
- Else choose highest scoring candidate gene found in cosmic census
- Else choose highest scoring protein coding candidate gene
- Else choose longest non-coding candidate gene

Finally, we filter the peaks to only highly significant deletions and amplifications using the following rules

- Deletions => Keep any peak with  $> 5$  homozygous deletions
- Amplifications => Keep any peak with score  $> 29$

These cut-offs were chosen using a binomial model which assumes the probability of any given gene being observed to be randomly deleted or highly amplified is equal to the average number of genes amplified or deleted in each event divided by the total number of genes considered. The cut-offs were chosen to be the lowest score with a q-value below 0.25. Since amplifications are generally much broader

(averaged genes affected per event of 41.6 compared to just 5.4 for deletions) a much higher number of genes is required to reach significance.

The calculation details for the cut-offs are presented in the table below.

	Cohort Data						Statistical Calculations				
	Count of Events	Sum Scores	Count of genes affected	Avg Genes affected per event	Avg score / event	Total Genes Tested	Probability event overlaps a given gene	Score cutoff	P Value of cutoff	Significant findings	Q Value
Dels	4,915	4,915	26,676	5.4	1.0	25,965	0.00021	5	0.00068	117	0.15
Amps	3,925	6,959	163,393	41.6	1.8	25,965	0.00160	29	0.00030	33	0.23

This model is likely to be highly conservative as it assumes that all the events are passengers, whereas in fact a high proportion contain driver genes.

## 16. Fragile site annotation

Homozygous deletions were also annotated as common fragile site (CFS) based on their genomic characteristics. This annotation is not definitive, but is useful as CFS are known to be regions of high genomic instability. Hence despite being significantly deleted, their status as a genuine cancer driver remains unclear.

There is no absolute agreement on which regions should be classified as CFS, but two well-known features are a strong enrichment in long genes and a high rate of observed deletions of up to 1 megabase<sup>18</sup>. Hence for this analysis we classified a gene as a fragile site if it met all the following criteria:

- Total length of gene > 500,000 bases
- More than 30% of all SVs with breakpoints that disrupt the gene are deletions with length greater than 20,000 bases and less than 1 megabase.
- The gene is not found to be significantly mutated (by dNdScv) in our cohort or in Martincorena et al.<sup>16</sup>

Using these criteria we annotated the following list of 16 Genes as fragile:

Gene	Chr	Start position	Length (bases)	Total Disruptive SV Count	% of SV that are DELs (>20kb & <1MB)
LRP1B	2	140,988,992	1,900,278	1,272	0.469
FHIT	3	59,735,036	1,502,097	2,128	0.596
LSAMP	3	115,521,235	2,194,860	1,306	0.364
NAALADL2	3	174,156,363	1,367,065	1,198	0.456
CCSER1	4	91,048,686	1,474,378	1,398	0.441
PDE4D	5	58,264,865	1,553,082	1,166	0.458
GMDS	6	1,624,041	621,885	399	0.441
PARK2	6	161,768,452	1,380,351	1,296	0.555
IMMP2L	7	110,303,110	899,463	1,028	0.444
PTPRD	9	8,314,246	2,298,477	1,264	0.309

PRKG1	10	52,750,945	1,307,165	781	0.318
GPHN	14	66,974,125	674,395	291	0.306
WWOX	16	78,133,310	1,113,254	1,319	0.541
MACROD2	20	13,976,015	2,057,827	3,039	0.605
DMD	X	31,115,794	2,241,764	789	0.328
DIAPH2	X	95,939,662	920,334	331	0.381

We also noted that 4 other significantly deleted genes (STS,HDHD1,LRRN3 and LINC00290), though not fulfilling the length criteria above have a particularly high proportion of deletion SVs between 20kb and 1 megabase (over 60%) and hence were also marked as fragile:

Gene	Chr	Start position	Length (bases)	Total Disruptive SV Count	% of SV that are DELs (>20kb & <1MB)
LINC00290	4	181,985,242	95,060	64	0.641
LRRN3	7	110,731,062	34,448	70	0.686
STS	X	7,137,497	135,354	168	0.649
HDHD1	X	6,966,961	99,270	126	0.659

Two of these genes (STS and HDHD1) fall in a previously identified CFS region (FRAXB) and a third, LRNN3, falls in another knowns CFS region (FRAX7). The final one, LINC00290 is a long non-coding RNA with an unknown status as cancer driver.

## 17. Driver catalog construction

We created a catalog of each and every driver in our cohort across all variant types on a per patient basis. This was done in a similar incremental manner to Sabarinathan et al<sup>19</sup> (N. Lopez, personal communication) whereby we first calculated the number of drivers in a broad panel of known and significantly mutated genes across the full cohort, and then assigned the drivers for each gene to individual patients by ranking and prioritising each of the observed variants. Key points of difference in this study were both the prioritisation mechanism used and our choice to ascribe each mutation a probability of being a driver rather than a binary cutoff based on absolute ranking.

The four detailed steps to create the catalog are described below:

### 1. Create a panel of driver genes for point mutations using significantly mutated genes and known drivers

We created a gene panel using the union of

- Martincorena significantly mutated genes<sup>16</sup> (filtered to significance of  $q<0.01$ )
- HMF significantly mutated genes ( $q<0.01$ ) at global level or at cancer type level
- Cosmic Curated Genes<sup>12</sup> (v83)

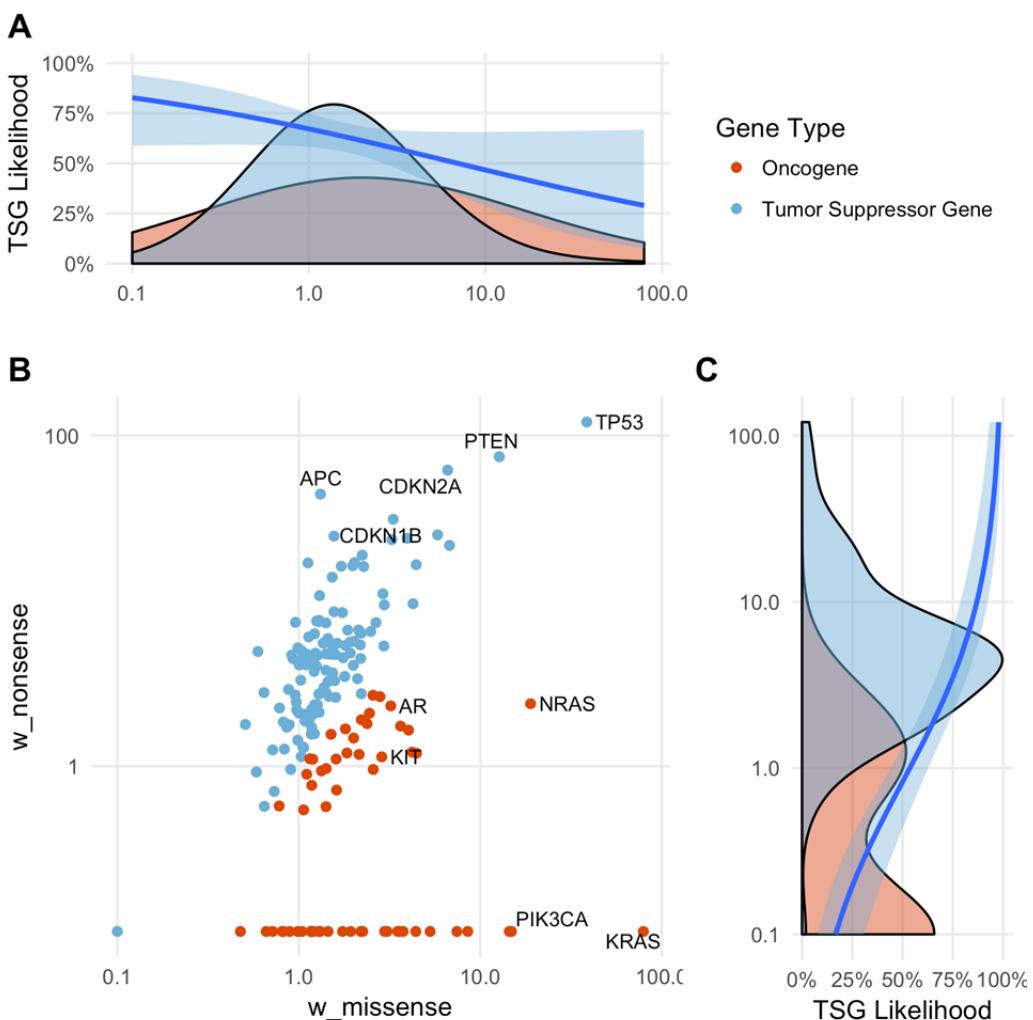
### 2. Determine TSG or Oncogene status of each significantly mutated gene

We used a logistic regression model to classify the genes in our pane as either tumor suppressor gene (TSG) or oncogene. We trained the model using unambiguous classifications from the Cosmic curated genes, i.e. a gene was considered either a Oncogene or TSG but not both. We determined that the dNdS missense and nonsense ratios ( $w_{\text{missense}}$  and  $w_{\text{nonsense}}$ ) are both significant predictors of the classification. The coefficients are given in the table below.

	Estimate	Std. Error	z value	Pr(> z )
intercept	0.1830	0.3926	0.466	0.64106
$w_{\text{missense}}$	-0.6869	0.2643	-2.599	0.00936
$w_{\text{nonsense}}$	0.5237	0.1116	4.691	2.72e-06

We applied the model to all significantly mutated genes in Matincorena and HMF as well as any ambiguous Cosmic curated genes.

The following figure shows all genes that have classified using the logistic regression model. Figures A and C show the likelihood of a gene being classified as a TSG under a single variate logistic model of  $w_{\text{missense}}$  and  $w_{\text{nonsense}}$  respectively. Figure B shows the classification after the multivariate regression using both predictors.



### 3. Add drivers from all variant classes to the catalog

Variants were added to the driver catalog which met any of the following criteria

- All missense and inframe indels for panel oncogenes
- All non synonymous and essential splice point mutations for tumor suppressor genes
- All high level amplifications (min exonic copy number > 3 \* sample ploidy) for both significantly amplified target genes and panel oncogenes
- All homozygous deletions for significantly deleted target genes and panel TSG (except for the Y chromosome as described before)
- All known or promiscuous inframe gene fusions as described above
- Recurrent TERT promoter mutations

### 4. Calculate a per sample driver likelihood for each gene in the catalog

A driver likelihood estimate between 0 and 1 was calculated for each variant in the gene panel to ensure that only excess mutations are used for determining the number of drivers in cancer cohort groups or at the individual sample level. High level amplifications, Deletions, Fusions, and TERT promoter mutations are all rare so were assumed to have a likelihood of 1 when found affecting a driver gene, but for coding mutations we need to account for the large number of passenger point mutations that are present throughout the genome and thus also in driver genes.

For coding mutations we also marked coding mutations that are highly likely to be drivers and/or highly unlikely to have occurred as passengers as driver likelihood of 1, specifically:

- Known hotspot variants
- Variants within 5 bases of a known pathogenic hotspot in oncogenes
- Inframe indels in oncogenes with repeat count < 8 repeats. Longer repeat count contexts are excluded as these are often mutated by chance in MSI samples
- Biallelic variants in tumor suppressor genes

For the remaining variants (non-hotspot missense variants in oncogenes and non-biallelic variants in TSG) these were only assigned a  $> 0$  driver likelihood where there was a remaining excess of unallocated drivers based on the calculated dNdS rates in that gene across the cohort after applying the above rules. Any remaining point mutations were assigned a driver likelihood between 0 and 1 using a bayesian statistic to calculate a sample specific likelihood of each gene based on the type of variant observed (missense, nonsense, splice or INDEL) and taking into account the mutational load of the sample. The principle behind the method is that the likelihood of a passenger variant occurring in a particular sample should be approximately proportional to the tumor mutational burden and hence variants in samples with lower mutational burden are more likely to be drivers.

The sample specific likelihood of a residual excess variant being a driver is estimated for each gene using the following formula:

$$P(\text{Driver}|\text{Variant}) = P(\text{Driver}) / (P(\text{Driver}) + P(\text{Variant}|\text{Non-Driver}) * (1-P(\text{Driver})))$$

where  $P(\text{Driver})$  in a given gene is assumed to be equal across all samples in the cohort, ie:

$$P(\text{Driver}) = (\text{residual unallocated drivers in gene}) / \# \text{ of samples in cohort}$$

And  $P(\text{Variant}|\text{Non-Driver})$ , the probability of observing  $n$  or more passenger variants of a particular variant type in a sample in a given gene, is assumed to vary according to tumor mutational burden, and is modelled as a poisson process:

$$P(\text{Variant}|\text{Non-Driver}) = 1 - \text{poisson}(\lambda = \text{TMB(Sample)} / \text{TMB(Cohort)} * (\# \text{ of passenger variants in cohort}), k=n-1)$$

All counts reported in the paper at a per cancer type or sample level refer to the sum of driver likelihoods for that cancer type or sample.

### 18. Driver co-occurrence analysis

To examine the co-occurrence of drivers, the driver-gene catalog was filtered to exclude fusions and any driver with a driver likelihood of < 0.5. Separately for each cancer type, every pair of driver genes was tested to see whether they co-occur more or less frequently than expected if they were independent using Fisher's Exact Test. The results were adjusted to a FDR using the number of gene-pair comparison being tested in each cancer type cohort. Gene pairs with a positive correlation which were on the same chromosome were excluded from the analysis as they are frequently co-amplified or deleted by chance.

### 19. Actionability analysis

To determine clinical actionability of the variants observed in each sample, we mapped all variants to 3 external clinical annotation databases

- OncoKB<sup>7</sup> (download = 01-mar-2018)
- CGI<sup>8</sup> (update: 17-jan-2018)
- CIViC<sup>5</sup> (download = 01-mar-2018)

In order to be able to aggregate and compare this data, we have mapped each of the databases to a common data model using the following rules:

#### 1. Level of evidence mapping

The 3 databases we used in this study define different level for evidence items, depending on evidence strength. In order to be able to aggregate and compare this data, we have mapped the CGI and OncoKB evidence levels on the CIViC evidence levels defined at: <https://civicdb.org/help/evidence/evidence-levels>.

HMF	CIViC	CGI	OncoKB
A	A	FDA guidelines, NCCN guidelines, NCCN/CAP guidelines, CPIC guidelines, European Leukemia Net guideline	1 2 R1
B	B	Clinical trials, Late trials, Late trials, Pre-clinical	3 R2
C	C	Early trials, Case report	
D	D	Pre-clinical	4, R3

In this study we considered only A and B level variants. This classification roughly corresponds to the recently proposed ESMO Scale for Clinical Actionability of molecular Targets (ESCAT)<sup>20</sup> in the following way:

HMF A: ESCAT I-A+B (for on label) and I-C (for off-label)

HMF B: ESCAT II-A+B (for on label) and III-A (for off-label)

## 2. Response type Mapping

We also mapped response type to a common data model. First we filtered out evidence items from the annotation databases that do not lead to clinical actionability (for example prognostic biomarkers). The remaining evidence items were mapped as either responsive or resistant based on the following rules:

HMF	CIViC	CGI	OncoKB
Responsive	Sensitivity	Responsive	1 2 3 4
Resistant	Resistant or Non-Response	Resistant	R1 R2 R3

## 3. Mutation/Event type mapping

Each evidence item was mapped to HMF data as one of 4 event types according to the following criteria

HMF Event type	Matching Criteria
Somatic Point Mutation	HGVS / genomic coordinates converted to chromosome, position, ref and alt and mapped to exact variants in our database
Somatic Range Event	Matched to missense / inframe variants in Oncogenes and any non-synonymous variant in TSG contained within a defined range, either exon level, transcript level or specific coordinates. Where a transcript was not specified, the canonical transcript was always used to map coordinates
Somatic CNA	'Deletion' mapped to homozygous deletions and 'Amplification' mapped to high level amplification (>3x sample ploidy)
Fusion	Exact matching to an inframe fusion in our database. For OncoKB 'loss-of-function' fusions were excluded

A small number of items from CIViC level B evidence level were deemed either not specific enough or insufficiently supportive of actionability for this study and were filtered:

- Evidence items supporting TP53, KRAS & PTEN as actionable
- Evidence items supporting actionability with 'chemotherapy' (ie. chemotherapy in general rather than a specific treatment), 'asprin' or 'steroids'

Finally, a number of suspicious fusions from each of the database were curated by either changing the 5' and 3' partners or filtered out altogether based on referring to the original evidence sources, specifically:

HMF Curation	CIViC	CGI	OncoKB
Filtered Fusions	BRAF - CUL1	RET - TPCN1	
5' and 3' partners exchanged		ABL1 - BCR PDGFRA - FIP1L1 PDGFB - COL1A1	ROS1 - CD74 EP300 - MLL EP300 - MOZ RET - CCDC6

Some of the more complex event types from the 3 databases have not been fully interpreted and have been excluded from this analysis.

#### 4. Cancer type mapping

Each evidence event mapped was also determined to be either on-label (ie. evidence supports treatment in that specific cancer type) or off-label (evidence exists in another cancer type) for each specific sample. To do this, we have annotated both the patient cancer types and the database cancer types with relevant DOIDs, using the disease ontology database available at: <http://disease-ontology.org>.

Patient cancer types from the HMF database were annotated according to the following table:

HMF tumor type	DOID
Biliary	4607
Bone/Soft tissue	201;9253
Breast	1612
CNS	3620;3070
Colon/Rectum	9256;219
CUP	-
Esophagus	5041;4944
Head and neck	11934;8618
Kidney	263;8411
Liver	3571
Lung	1324
Mesothelioma	1790
NET	-
Other	-
Ovary	2394
Pancreas	1793
Prostate	10283
Skin	4159
Stomach	10534
Urinary tract	3996
Uterus	363

Database cancer types were mapped to a DOID by automatically querying the ontology on the disease names. Some CIViC evidence items are already annotated with a DOID in the database, this was used if present. We also manually annotated with DOIDs some of the database cancer types that failed the automatic query:

cancerType	DOID	Ontology term
All Tumors	162	cancer
Any cancer type	162	cancer
B cell lymphoma	707	B-cell lymphoma
Billiary tract	4607	biliary tract cancer
Bladder	11054	urinary bladder cancer
Cervix	4362	cervical cancer
CNS Cancer	3620	central nervous system cancer
Dedifferentiated Liposarcoma	3382	liposarcoma
Endometrium	1380	endometrial cancer
Esophagogastric Cancer	5041	esophageal cancer
Gastrointestinal stromal	9253	gastrointestinal stromal tumor
Giant cell astrocytoma	3069	astrocytoma

Hairy-Cell leukemia	285	hairy cell leukemia
Head and neck	11934	head and neck cancer
Head and neck squamous	5520	head and neck squamous cell carcinoma
Hepatic carcinoma	686	liver carcinoma
Hepatocellular Mixed Fibrolamellar Carcinoma	0080182	mixed fibrolamellar hepatocellular carcinoma
Inflammatory myofibroblastic	0050905	inflammatory myofibroblastic tumor
Lung	1324	lung cancer
Lung squamous cell	3907	lung squamous cell carcinoma
Melanoma	8923	Skin melanoma
Mesothelioma	1790	malignant mesothelioma
Neuroendocrine	169	neuroendocrine tumor
Non-small cell lung	3908	non-small cell lung carcinoma
Ovary	2394	ovarian cancer
Pancreas	1793	pancreatic cancer
Renal	263	kidney cancer
Salivary glands	8850	salivary gland cancer
Stomach	10534	stomach cancer
Thymic	3277	thymus cancer
Thyroid	1781	thyroid cancer
Well-Differentiated Liposarcoma	3382	liposarcoma

In case a matching DOID was found for the disease, we annotated the disease with a DOID set consisting of: the disease DOID, all the children DOIDs and all the parent disease DOIDs.

A treatment is defined as on-label if any of the DOIDs of the patient cancer is present in the DOID set of the disease.

## 5. MSI actionability

Samples classified as MSI in our driver catalog were also mapped as actionable at level A evidence based on clinical annotation in the OncoKb database

## 6. Aggregation of evidence

For each actionable mutation in each sample, we aggregated all the mapped evidence that was available supporting both on-label and off-label treatments at an A or B evidence level. Treatments that also had evidence supporting resistance based on other biomarkers in the sample at the same or higher level were excluded as non-actionable.

For each sample we reported the highest level of actionability, ranked first by evidence level and then by on-label vs off-label.

## 20. Data availability

All data described in this study is freely available for academic use from the Hartwig Medical Foundation through standardized procedures and request forms which can be found at

<https://www.hartwigmedicalfoundation.nl>. Briefly, a data request can be initiated by filling out the standard form in which intended use of the requested data is motivated. First, an advice on scientific feasibility and validity is obtained from experts in the field which is used as input by an independent Data Access Board who also evaluates if the intended use of the data is compatible with the consent given by the patients and if there would be any applicable legal or ethical constraints. Upon formal approval by the Data Access Board, a standard license agreement which does not have any restrictions regarding Intellectual Property resulting from the data analysis needs to be signed by an official organisation representative

before access to the data is granted. Raw data files will be made available through a dedicated download portal with two-factor authentication. As of August 2018, more than 30 national and international requests have already been granted through this mechanism.

## 21. References

1. Bins, S. *et al.* Implementation of a Multicenter Biobanking Collaboration for Next-Generation Sequencing-Based Biomarker Discovery Based on Fresh Frozen Pretreatment Tumor Tissue Biopsies. *Oncologist* **22**, 33–40 (2017).
2. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
3. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinformatics* **43**, 11.10.1–33 (2013).
4. Saunders, C. T. *et al.* Strelka: accurate somatic small-variant calling from sequenced tumor-normal sample pairs. *Bioinformatics* **28**, 1811–1817 (2012).
5. Griffith, M. *et al.* CIViC is a community knowledgebase for expert crowdsourcing the clinical interpretation of variants in cancer. *Nat. Genet.* **49**, 170–174 (2017).
6. Tamborero, D. *et al.* Cancer Genome Interpreter annotates the biological and clinical relevance of tumor alterations. *Genome Med.* **10**, 25 (2018).
7. Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol* **2017**, (2017).
8. Cleveland, M. H., Zook, J. M., Salit, M. & Vallone, P. M. Determining Performance Metrics for Targeted Next-Generation Sequencing Panels Using Reference Materials. *J. Mol. Diagn.* (2018). doi:10.1016/j.jmoldx.2018.04.005
9. Eijkelenboom, A. *et al.* Reliable Next-Generation Sequencing of Formalin-Fixed, Paraffin-Embedded Tissue Using Single Molecule Tags. *J. Mol. Diagn.* **18**, 851–863 (2016).
10. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* **32**, 1220–1222 (2016).
11. Zerbino, D. R. *et al.* Ensembl 2018. *Nucleic Acids Res.* **46**, D754–D761 (2018).
12. Forbes, S. A. *et al.* COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids Res.* **45**, D777–D783 (2017).
13. Nilsen, G. *et al.* Copynumber: Efficient algorithms for single- and multi-track copy number segmentation. *BMC Genomics* **13**, 591 (2012).
14. Huang, M. N. *et al.* MSIseq: Software for Assessing Microsatellite Instability from Catalogs of Somatic Mutations. *Sci. Rep.* **5**, 13321 (2015).
15. Al-Kaabi, A., van Bockel, L. W., Pothen, A. J. & Willems, S. M. p16INK4A and p14ARF gene promoter hypermethylation as prognostic biomarker in oral and oropharyngeal squamous cell carcinoma: a review. *Dis. Markers* **2014**, 260549 (2014).
16. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**, 1029–1041 e21 (2017).
17. Mermel, C. H. *et al.* GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome Biol.* **12**, R41 (2011).
18. Glover, T. W., Wilson, T. E. & Arlt, M. F. Fragile sites in cancer: more than meets the eye. *Nat. Rev. Cancer* **17**, 489–501 (2017).
19. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *BioArchive* (2017). doi:10.1101/190330
20. Mateo, J. *et al.* A framework to rank genomic alterations as targets for cancer precision medicine: the ESMO Scale for Clinical Actionability of molecular Targets (ESCAT). *Ann. Oncol.* (2018).

Priestley, Baber, et al.  
Pan-cancer whole genome analyses of metastatic solid tumors  
doi:10.1093/annonc/mdy263