

1 Sex Differences in Oncogenic Mutational 2 Processes

3

4 Constance H. Li^{1,2}, Stephenie D. Prokopec¹, Ren X. Sun^{1,3}, Fouad Yousif¹, Nathaniel Schmitz¹,
5 Paul C. Boutros^{1,2,3,4,5,6,7,8}, for the PCAWG Molecular Subtypes and Clinical Correlates Working
6 Group, ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network

7

8 ¹ Ontario Institute for Cancer Research, Toronto, Ontario, Canada

9 ² Department of Medical Biophysics, University of Toronto; Toronto, Ontario, Canada

10 ³ Department of Pharmacology & Toxicology, University of Toronto; Toronto, Ontario, Canada

11 ⁴Vector Institute for Artificial Intelligence, Toronto, Canada

12 ⁵ Department of Human Genetics, University of California, Los Angeles, CA, USA

13 ⁶ Department of Urology, University of California, Los Angeles, CA, USA

14 ⁷ Jonsson Comprehensive Cancer Center, University of California, Los Angeles, CA, USA

15 ⁸ Institute for Precision Health, University of California, Los Angeles, CA, USA

16

17 Correspondence should be addressed to P.C.B. (PBoutros@mednet.ucla.edu)

18 12-109 CHS

19 10833 Le Conte Avenue

20 Los Angeles, CA, 90095

21 Phone: 310-794-7160

22

23

24 Abstract

25 Sex differences have been observed in multiple facets of cancer epidemiology, treatment and
26 biology, and in most cancers outside the sex organs. Efforts to link these clinical differences to
27 specific molecular features have focused on somatic mutations within the coding regions of the
28 genome. Here, we describe the first pan-cancer analysis of sex differences in whole genomes of
29 1,983 tumours of 28 subtypes from the ICGC Pan-Cancer Analysis of Whole Genomes project.
30 We both confirm the results of exome studies, and also uncover previously undescribed sex
31 differences. These include sex-biases in coding and non-coding cancer drivers, mutation
32 prevalence and strikingly, in mutational signatures related to underlying mutational processes.
33 These results underline the pervasiveness of molecular sex differences and strengthen the call
34 for increased consideration of sex in cancer research.

35 Sex disparities in cancer epidemiology include an increased overall cancer risk in males
36 corresponding with higher incidence in most tumor types, even after adjusting for known risk
37 factors^{1,2}. Cancer mortality is also higher in males, due in part to better survival for female patients
38 in many cancer types, including those of the colon and head & neck³. Interestingly, female
39 colorectal cancer patients respond better to surgery⁴ and adjuvant chemotherapy, though this is
40 partially due to biases in tumour location and microsatellite instability⁵. Similarly, premenopausal
41 female nasopharyngeal cancer patients have improved survival regardless of tumour stage,
42 radiation or chemotherapy regimen⁶. There is a growing body of evidence for sex differences in
43 cancer genomics⁷⁻¹³, but their molecular origins and clinical implications remain largely elusive.

44 Previous studies have mostly focused on protein coding regions, leaving the vast majority of the
45 genome unexplored. We hypothesized that there are uncharacterized sex differences in the non-
46 coding regions of the genome. Using whole genome sequencing data from the Pan-cancer
47 Analysis of Whole Genomes (PCAWG) project¹⁴, we performed a survey of sex-biased mutations
48 in 1,983 samples (1,213 male, 770 female) from 28 tumour subtypes, excluding those of the sex
49 organs (**Supplementary Table 1**). We also excluded the X and Y chromosomes to focus on
50 autosomal sex differences in cancers affecting both men and women, but there are known to be
51 significant X-chromosome mutational differences between tumours arising in men and women¹⁵.
52 Our analysis revealed sex differences in both genome-wide phenomena and in specific genes.
53 These sex-biases occur not only at the pan-cancer level across all 1,983 samples, but also in
54 individual tumour subtypes.

55 **Sex-biases in driver genes, mutation load and tumour evolution**

56 We began by investigating sex differences in driver gene mutation frequencies, focusing on 165
57 coding and nine non-coding mutation events¹⁶ (**Supplementary Table 2**). We used proportions
58 tests to identify candidate sex-biased events with a false discovery rate (FDR) threshold of 10%.
59 These putative sex-biased events were then modeled using logistic regression (LGR) to adjust
60 for tumour subtype, ancestry and age (**Online Methods**). We found several sex-biased pan-
61 cancer driver events, including *CTNNB1* which was mutated in 5.0% more male-derived than
62 female-derived tumours (male: 7.6%, female: 2.7%, 95% CI: 2.9-7.0%, prop-test $q = 3.6 \times 10^{-4}$,
63 LGR $q = 5.0 \times 10^{-3}$; **Figure 1a, left**). *ALB* was also mutated in a larger proportion of male-derived
64 tumours (male: 3.2%, female: 0.54%, 95% CI: 1.4-3.9%, prop-test $q = 0.0038$, LGR $q = 6.5 \times 10^{-3}$)
65 while in contrast, *PTCH1* (male: 0.44%, female: 2.0%, 95% CI: 0.38-2.8%, prop-test $q = 0.028$,
66 LGR $q = 0.011$) was mutated in more female-derived samples.

67 We also identified tumour subtype-specific sex-biased driver mutations (**Figure 1a, right**).
68 Similarly to the pan-cancer driver analysis, we first identified putative sex-biases using proportions
69 tests and a 10% FDR threshold, and followed up with tumour subtype-specific logistic regression
70 models (model descriptions in **Supplementary Table 1**). *CTNNB1* mutation frequency was sex-
71 biased in liver hepatocellular cancer (Liver-HCC), again with more male-derived samples
72 harbouring *CTNNB1* mutations: (male: 31%, female: 13%, 95% CI: 8.1-28%, prop-test $q = 0.047$,
73 LGR $q = 8.2 \times 10^{-3}$, **Figure 1a, right**). This mirrors our previous finding of sex-biased *CTNNB1*
74 mutation frequency in liver cancer from TCGA exome sequencing data, with similar effect sizes
75 (male: 33% vs. female: 12%¹¹). The largest sex-disparity was in a non-coding driver event in
76 thyroid cancer (Thy-AdenoCA): *TERT* promoter mutations were observed in 64% of male-derived
77 samples compared with only 11% of female-derived samples (95% CI: 18-89%, prop-test $q =$

78 6.9×10^{-3} , LGR $q = 0.074$, **Figure 1a, right**), again supporting a previous finding¹⁷. Other putative
79 sex-biased events were detected, but were not statistically significant after multivariate
80 adjustment at present sample-sizes (**Supplementary Table 2**). These results demonstrate that
81 mutation of key cancer-driving genes is sex-biased both within and across specific tumour
82 subtypes.

83 Our previous work¹² found sex biased mutation density across a number of tumour subtypes,
84 including cancers of the liver, kidney and skin. We therefore investigated mutation density here
85 to identify tumour subtypes where the cancer genomes of one sex accumulates more somatic
86 single nucleotide variants (SNVs) than those of the other sex, and whether these sex-biases might
87 be related to sex-biased driver gene mutation frequency. Returning to our statistical framework,
88 we first used univariate tests to identify putative sex-biases, and then applied multivariate linear
89 regression (LNR) on Box-cox transformed mutation load to adjust for possible confounders. The
90 Box-cox transformation applies a power function to modify the shape of a variable's distribution
91 to better approximate a normal distribution (**Online Methods**). We also compared the total
92 number of somatic SNVs and further divided mutations by coding and non-coding SNVs to
93 determine whether sex-biases may be influenced by specific genomic contexts. Across all pan-
94 cancer samples, we found higher mutation prevalence in male-derived samples in all three
95 contexts (coding: difference in location = 0.41 mut/Mbp, 95% CI = 0.28-0.54 mut/Mbp, u-test $q =$
96 2.2×10^{-10} , LNR $q = 7.5 \times 10^{-4}$; non-coding: difference in location = 0.60 mut/Mbp, 95%CI = 0.43-
97 0.80 mut/Mbp, u-test $q = 7.9 \times 10^{-11}$, LNR $q = 6.5 \times 10^{-4}$; overall: difference in location = 0.60
98 mut/Mbp, 95%CI = 0.42-0.79 mut/Mbp, u-test $q = 7.5 \times 10^{-11}$, LNR $q = 1.9 \times 10^{-6}$; **Supplementary**
99 **Table 3**). These sex-biases remained significant even after adjusting for tumour subtype, ancestry
100 and age in multivariate analysis (**Figure 1b, left**), demonstrating robust sex-biases in pan-cancer
101 mutation prevalence across different contexts.

102 We also investigated somatic SNV burden in each of the 23 individual tumour subtypes with at
103 least 15 samples, applying the same statistical approach with tumour subtype-specific models
104 (model descriptions in **Supplementary Table 1**). We found sex-biased mutation load in three
105 tumour subtypes (**Figure 1b, right**), with higher male coding mutation load in thyroid cancer
106 (difference in location = 0.26 mut/Mbp, 95%CI = 0.12-0.43 mut/Mbp, u-test $q = 0.028$, LNR $q =$
107 0.041), and higher male load in hepatocellular cancer and kidney renal cell cancer (Kidney-RCC)
108 in all three contexts (**Supplementary Table 3**). We compared the group rank differences of
109 coding and non-coding mutation load between the sexes and found that in renal cell cancer, the
110 differences were similar at 0.40 mut/Mbp for non-coding mutations and 0.37 mut/Mbp for coding
111 mutations. In hepatocellular cancer however, the median sex-difference in non-coding mutation
112 load was higher than the difference in coding mutation load (non-coding difference = 0.84
113 mut/Mbp vs. coding difference = 0.53 mut/Mbp). There is a similar effect in pan-cancer mutation
114 load (non-coding difference = 0.60 mut/Mbp vs coding difference = 0.41 mut/Mbp) suggesting
115 mutation context may play a role in sex-biased SNVs in some tumour subtypes.

116 To determine whether sex-biased mutation load may be associated with sex-biased driver gene
117 mutation frequency, we focused on each driver gene and investigated SNV burden in the relevant
118 tumour subtype. We did not find significant relationships between SNV burden and mutations in
119 *PTCH1*, *ALB*, *CTNNB1* in pan-cancer analysis, nor was there an association for *CTNNB1*
120 mutation in hepatocellular cancer. In thyroid cancer however, *TERT* promoter mutation was

121 associated with increased coding mutation burden (median_{TERT-wt} = 0.26 mut/Mbp vs median_{TERT-}
122 _{mut} = 0.66 mut/Mbp, u-test p = 4.9x10⁻⁶). We used a linear regression model to determine if the
123 sex-bias in coding mutation load could be explained by *TERT* mutation frequency and found this
124 was indeed the case (linear regression p_{TERT} = 2.4x10⁻⁵, p_{sex} = 0.37, **Figure 1c**). In addition, we
125 examined matched mutation timing data and found that of eleven samples with *TERT* promoter
126 mutations, nine of these were truncal events, suggesting that an early sex-bias in *TERT* promoter
127 mutation frequency is associated with sex-biased coding mutation load in this tumour subtype.
128 Indeed, mutations in all sex-biased driver genes were overwhelmingly truncal events.

129 We then asked if these driver mutations might occur at different stages of tumour evolution
130 between men and women, and started with tumour evolution structure. We compared the
131 proportions of polyclonal vs. monoclonal tumours between the sexes and did not find significant
132 sex differences in the proportions of polyclonal tumours bearing mutations in *PTCH1*, *ALB* or
133 *CTNNB1* for sex-biased pan-cancer drivers, or in *TERT* promoter-mutated samples in thyroid
134 cancer (**Supplementary Figure 1**). We did detect a putative bias in the proportion of polyclonal
135 *CTNNB1*-mutated samples in hepatocellular cancer (80% of male-derived samples are polyclonal
136 vs. 46% of female-derived samples, 95%CI = -0.019 – 0.70, prop-test p = 0.039), and accounted
137 for polyclonality when comparing the timings of the mutations in these driver events. On
138 subsequently examining the frequency of clonal vs. subclonal driver mutation events between the
139 sexes, we found that while there were differences in the proportions of truncal mutations (eg:
140 100% of *TERT* promoter mutations were truncal events in male-derived vs. 50% truncal events in
141 female-derived thyroid cancer patients), no comparisons were statistically significant.

142 Broadening beyond sex-biased driver mutations, we expanded our clonality analysis to perform
143 a general survey of clonal structure and mutation timing across all tumour subtypes and mutations
144 (**Supplementary Table 4**). We found that female-derived biliary adenocarcinoma (Biliary-
145 AdenoCA) tumours were frequently polyclonal, while most male-derived tumours were
146 monoclonal (26% male-derived samples are polyclonal vs. 80% female-derived, 95% CI = 19 –
147 88%, prop-test q = 0.063, LGR q = 0.024; **Figure 1d**). In addition, we found intriguing evidence
148 suggesting there may be sex-differences in the mutation timing of structural variants in this tumour
149 subtype. Structural variants (SVs) in male-derived samples tended to be truncal events more
150 frequently than in female-derived samples (median male percent truncal SVs = 100% vs. median
151 female = 82%, u-test q = 0.081, LNR q = 0.024; **Figure 1e**). Though other comparisons did not
152 reach our statistical significance threshold, we found some interesting trends that may merit future
153 study, including in esophageal cancer (Eso-AdenoCA) where SVs in female-derived samples
154 were more frequently truncal events while SVs in male-derived samples occurred more frequently
155 in subclones (median male percent truncal SVs = 55%, median female = 100%; **Supplementary**
156 **Figure 2**), and in medulloblastoma, where insertion-deletions (indels) were more frequently
157 truncal events in female-derived samples than male (median male percent of truncal indels =
158 65%, median female proportion of truncal indels = 70%; **Supplementary Figure 3**). Our analysis
159 of sex differences in tumour evolution identified some sex-biased events and also hint at putative
160 sex-biases that should be further explored in future analyses.

161 **Sex-biases in genome instability and CNAs**

162 Next, we examined percent genome altered (PGA), which provides a summary of copy number
163 aberration (CNA) load. A proxy for genome instability, PGA is a complementary measure of

164 mutation density to somatic SNV burden. While we did not find associations between sex and
165 autosome-wide PGA, we observed sex-biases in the copy number burden for specific
166 chromosomes (**Figure 2a**). In pan-cancer analysis, male-derived samples exhibited a slight but
167 significant higher percent chromosome altered for chromosome 7 even after accounting for
168 tumour subtype, ancestry and age using linear regression (median male PGA-7 = 5.4%, median
169 female PGA-7 = 0.37%, difference in location = 0.0037%, 95%CI = 9.4×10^{-4} – 2.4×10^{-3} %, u-test =
170 5.0×10^{-3} , LNR q = 0.027; **Supplementary Table 5**). In individual tumour subtypes, we found sex-
171 biased PGA in renal cell cancer (chromosomes 7 & 12) and hepatocellular cancer (chromosomes
172 1 & 16). By looking at copy number gains and losses separately, we additionally identified
173 chromosomes with sex-biases in the burden of copy number gains and losses (**Supplementary**
174 **Figure 4, Supplementary Table 5**), including sex-biased percent copy gained on chromosomes
175 5, 8 and 17 in pan-cancer samples.

176 We next compared CNA frequency on the gene level to identify genes lost or gained at sex-biased
177 rates. Across all pan-cancer samples, we found 4,285 sex-biased genes across 15 chromosomes
178 (**Figure 2b, Supplementary Tables 6 & 7**, LGR q-value < 10%). These genes were all more
179 frequently gained in male-derived samples than female with a difference in copy number gain
180 frequency reaching ~10% on chromosomes 7 and 8. Genes with male-dominated copy number
181 gains include the oncogenes *MYC* (male gain frequency = 37% vs. female gain frequency = 28%,
182 95% CI = 5.2-14%, prop-test q = 2.5×10^{-3} , LGR q = 0.068) and *ERBB2* (male gain frequency =
183 21% vs. female gain frequency = 16%, 4.7%, 95% CI = 1.1-8.3%, prop-test q = 0.041, LGR q =
184 0.088). The driver *CTNNB1* was also more frequently gained in male samples (male gain
185 frequency = 8.9% vs. female gain frequency = 5.2%, 95% CI = 1.4-6.1%, prop-test q = 0.016,
186 LGR q = 0.053), mirroring our finding of higher male pan-cancer mutation frequency on the SNV
187 level for this oncogene. We did not find pan-cancer sex-biased copy number losses.

188 We repeated this analysis for every tumour subtype independently and found sex-biased CNAs
189 in renal cell and hepatocellular cancer (**Supplementary Tables 6 & 7**). In renal cell cancer, the
190 1,986 sex-biased gains all occurred more frequently in male-derived samples, with differences in
191 frequency up to 35% (**Figure 2c**). They spanned across chromosomes 7 and 12, agreeing with
192 our finding of male-dominated genome instability in these chromosomes (**Figure 2a,**
193 **Supplementary Figure 4**). In contrast to the male-dominated gain pan-cancer and renal cell
194 findings, we found higher female frequency of copy number losses in hepatocellular cancer
195 (**Figure 2d**). We identified 2,610 genes with higher copy number loss rates in female-derived
196 samples. As observed in renal cell cancer, some of these losses span whole chromosomes, in
197 this case chromosomes 3 and 16. Other sex-biased losses were found only across one
198 chromosome arm (1p, 4q) or as focal events (eg. *PCDH9* on chromosome 13). The sex-biased
199 gene-level events on chromosomes 1 and 16 agreed with the sex-biased genome instability
200 findings but on returning to the PGA analysis, we found that chromosomes 3 and 4 had trending
201 sex-biased genome instability (u-test q < 0.2, **Figure 1a, Supplementary Table 5**), suggesting
202 that sex-biased PGA may guide identification of sex-biased CNAs on the gene level.

203 Thus, using sex-biased PGA as a guide, we more closely examined regions of interest in tumour
204 subtypes of that did not have sex-biased CNAs in our general CNA analysis, but did have possible
205 sex-biased genome instability (u-test q < 0.2): biliary cancer, B-cell non-Hodgkin lymphoma
206 (Lymph-BNHL), chronic lymphocytic leukemia (Lymph-CLL) and melanoma (Skin-Melanoma).

207 We found an additional 203 genes on the p-arm of chromosome 8 that were more frequently lost
208 in female-derived samples in biliary cancer (**Supplementary Figure 5**). These copy number
209 losses were 50% more common in female-derived samples and affect genes such as *DLC1*, a
210 known tumour suppressor in hepatocellular cancer that is thought to play a similar role in
211 gallbladder cancer¹⁸. While we did not identify additional sex-biased CNAs in non-Hodgkin
212 lymphoma, chronic lymphocytic leukemia or melanoma, the sex-biased PGA results suggest
213 these as regions of interest for future work. Thus in addition to sex-biased SNV events, we also
214 identified sex-biased CNAs from this whole genome sequencing data.

215 **Sex biases in mutation signatures**

216 We hypothesized that sex differences in mutation load and tumour evolution characteristics may
217 be driven by varying mutational processes. In addition to single base substitution (SBS)
218 signatures, which have been well annotated and linked to tumour aetiology^{19,20}, we also examined
219 doublet base substitution (DBS) and small insertion-deletion (ID) signatures. Sex differences in a
220 mutational signature could shine insight on molecular differences between the sexes. For each of
221 47 validated PCAWG SBS, 11 DBS, and 17 ID signatures²¹, we performed a two-stage analysis.
222 We first compared the proportions of signature-positive samples between the sexes; that is, we
223 looked at the proportions of samples with any mutations attributed to the signature to determine
224 whether there was a relationship between each signature and sex. Then, we focused on
225 signature-positive samples and compared the percentage of mutations attributed to each
226 signature between the sexes. For both analyses, we used univariate techniques to identify
227 putative events and adjusted for additional variables using linear models.

228 At the pan-cancer level, we found eight signatures that occurred more frequently in one sex over
229 the other (**Figure 3a, Supplementary Table 8**). In particular, SBS1 was more common in female-
230 derived samples (89% of male-derived vs. 97% of female-derived, χ^2 -test $q = 3.9 \times 10^{-10}$, LGR $q =$
231 5.1×10^{-7}) and was also associated with a higher percentage of mutations in these samples (male
232 median percent mutations attributed to SBS1 = 8.4%, female median = 10%, u-test $q = 0.026$,
233 LNR $q = 0.021$). SBS1 is thought to be caused by deamination of 5-methylcytosine to thymine,
234 resulting in base substitutions. Though it is correlated with age, our multivariate model accounts
235 for this variable and the sex-bias remains even adjusting for age. SBS40 was also detected in a
236 larger proportion of female-derived samples (42% of male-derived vs. 52% of female-derived, χ^2 -
237 test $q = 1.7 \times 10^{-4}$, LNR $q = 0.08$), though we did not find a difference in the percentage of attributed
238 mutations (u-test $q = 0.17$). Other sex-biased SBS signatures include SBS16, SBS17a and
239 SBS17b, which were all more frequently detected in male-derived samples, and SBS40, which
240 was more frequent in female-derived samples. These signatures are of unknown aetiology.

241 One ID signature was detected at different rates between the sexes, and two ID signatures had
242 different rates of attributed mutations. ID8 occurred more frequently in male-derived samples
243 (53% of male-derived vs. 47% of female-derived, χ^2 -test $q = 0.068$, LGR $q = 0.018$) though there
244 was no difference in the percentage of mutations attributed to either signature. The aetiology
245 underlying ID8 is not known, but this signature is thought to be associated with double strand
246 break repair where ID8-associated mutations resemble those related to radiation-induced
247 damage. Conversely, ID1 and ID5 were detected at similar frequencies between the sexes, but
248 had higher percentages of attributed mutations in female-derived samples. Mutations associated

249 with ID1 are thought to result from slippage during DNA replication and are associated with
250 defective DNA mismatch repair, suggesting that while male- and female-derived tumours harbour
251 defective DNA repair at similar rates, it is responsible for a larger proportion of mutations in
252 female-derived tumours.

253 Since mutational processes are disease-specific, we repeated the mutational signatures analysis
254 in each tumour subtype, again by first using univariate techniques to find putatively sex-biased
255 signatures, and then using linear models to adjust for age and ancestry. We identified six sex-
256 biased signatures in hepatocellular cancer (**Figure 3b, Supplementary Table 8**). Similar to our
257 pan-cancer finding, we again detected female-dominated bias in the proportion of SBS1-positive
258 samples (58% of male-derived vs. 88% of female-derived, χ^2 -test $q = 3.5 \times 10^{-5}$, LGR $q = 9.2 \times 10^{-6}$)
259 male-dominated bias in and SBS16 (16% of male-derived vs. 2.2% of female-derived, χ^2 $q =$
260 9.6×10^{-3} , LGR $q = 6.4 \times 10^{-3}$). There were four sex-biased ID signatures in this tumour subtype: ID3
261 (94% of male-derived vs. 81% of female-derived, χ^2 -test $q = 5.2 \times 10^{-3}$, LGR $q = 3.8 \times 10^{-3}$), ID8
262 (93% of male-derived vs. 78.7% of female-derived, χ^2 -test $q = 3.7 \times 10^{-3}$, LGR $q = 3.8 \times 10^{-3}$) and
263 ID11 (17% of male-derived vs. 1.1% of female-derived, χ^2 -test $q = 3.7 \times 10^{-3}$, LGR $q = 6.4 \times 10^{-3}$)
264 occurred more frequently in male-derived samples. While ID1 was detected at similar rates
265 between the sexes, a greater proportion of ID1-attributed mutations were found in female-derived
266 than male-derived samples (male median percent mutations attributed to ID1 = 21%, female
267 median = 27%, u-test $q = 2.0 \times 10^{-5}$, LR $q = 2.2 \times 10^{-6}$). As previously described, SBS1 and ID1 are
268 associated with base deamination and defective DNA mismatch repair. ID3 is associated with
269 tobacco smoke, and ID8 with double-stranded break repair. Taken together, sex-biases in the
270 aetiology underlying the molecular landscape of hepatocellular cancer begin to emerge. In this
271 tumour subtype, spontaneous or enzymatic deamination of 5-methylcytosine to thymine and
272 defective mismatch repair occur more frequently in female patients and are also responsible for
273 more mutations. Conversely, tobacco smoking is more common in male patients though the
274 number of mutations attributed to tobacco smoke is not different between the sexes; this leads to
275 more tobacco-associated male hepatocellular tumours.

276 In B-cell non-Hodgkin lymphoma, we identified a significant difference in the proportion of samples
277 with SBS17b-attributed mutations (**Figure 3c, Supplementary Table 8**). More male-derived
278 samples had mutations associated with this signature of unknown aetiology (57% of male-derived
279 vs. 25% of female-derived, χ^2 -test $q = 0.051$, LGR $q = 6.3 \times 10^{-4}$). There were also several intriguing
280 sex-differences in mutational signatures that did not meet our significance threshold. For instance,
281 DBS2 accounts for a higher percentage of mutations in male-derived samples (male median
282 percent mutations attributed to DBS2 = 50%, female median = 33%, **Supplementary Table 8**).
283 DBS2's association with tobacco smoking suggests that future insight in this signature may
284 provide molecular explanations for the sex-specific associations between smoking and thyroid
285 cancer risk²². As the aetiologies of these mutational signatures become better known, we can
286 better approach the causes of molecular sex differences underlying cancer aetiology and
287 progression. In particular, we may be able to discern environmental and lifestyle factors even in
288 the absence of reported data, and connect known risk factors with newly described mutational
289 processes.

290 Finally, to ensure that our findings were not skewed by differences in sequencing quality, we
291 checked for sex-biases in quality control (QC) metrics. These included comparing the coverage,

292 percentage of paired reads mapping to different chromosomes, and overall quality summary of
293 both tumour and normal genomes. We mirrored our main analyses and used u-tests or χ^2 tests
294 and linear modeling to check each QC metric. We did not find sex-biases in any QC metric in pan-
295 cancer or tumour subtype analysis after multiple adjustment except in raw somatic mutation
296 calling (SMC) coverage. SMC coverage was higher in male-derived samples in six tumour
297 subtypes including thyroid cancer and esophageal cancer, and was higher in female-derived
298 samples in lung adenocarcinoma and B-cell non-Hodgkin lymphoma (**Supplementary Table 9**).
299 While we do not find sex differences in comparing the SMC coverage pass/fail rates using a
300 recommended minimum of 2.6 gigabases covered, it is prudent to consider sex-biased SMC in
301 relation to our findings. There are also no sex-differences in the proportions of samples passing
302 quality checks for any other QC metric (**Supplementary Table 9, Supplementary Figures 6**).

303 Our analysis of whole genome sequencing data from the PCAWG project uncovered sex
304 differences in the largely unexplored non-coding autosomal genome. We found these biases in
305 measures of mutational load, tumour evolution, mutational signatures, and at the gene level.
306 While the majority of our findings describe pan-cancer differences, we have also uncovered an
307 intriguing glimpse into tumour subtype-specific differences. These tumour subtype-specific results
308 are limited by subtype sample size, and limited available annotation restricts the ability to account
309 for confounding variables. It is important to consider these results in context of the multivariable
310 models used, which do not directly capture characteristics such as tobacco smoking history or
311 tumour stage at diagnosis. Future increases in sample size and robust associated annotation will
312 allow for the detection of smaller effects and the control of more confounders. Nevertheless, our
313 analyses of driver genes and copy number alterations suggest functional impacts of genomic sex-
314 biases on the transcriptome and tumorigenesis. By using signatures to distinguish between
315 mutations attributed to lifestyle factors such as smoking, we can better describe sex differences
316 related to biological factors such as hormone activity. And despite low tumour subtype-specific
317 sample numbers, our mutation timing and mutational signatures findings at both the pan-cancer
318 and tumour-subtype level hint at underlying mutational processes that may give rise to molecular
319 sex-biases. Combined with our previous work in whole exome sequencing, we present a
320 landscape of sex-biases in cancer genomics and mutational processes (**Figure 4,**
321 **Supplementary Figure 7**).

322 It is becoming clear that sex differences occur across many mutation classes and the portrait of
323 differences for each tumour subtype is a unique reflection of active mutational processes and
324 tumour evolution. We have performed here the first pan-cancer analysis of sex differences in
325 whole genome sequencing data and catalogued previously undescribed sex-biases. However,
326 increased study of molecular sex differences in future large-scale sequencing efforts is needed
327 to strengthen the findings we present here, to determine why men and women have molecularly
328 different tumours, and to determine how this information can be leveraged to improve patient care.

329 Acknowledgments

330 The authors thank all the members of the Boutros lab for insightful discussions. This study was
331 conducted with the support of the Ontario Institute for Cancer Research to P.C.B. through funding
332 provided by the Government of Ontario. This work was supported by the Discovery Frontiers:
333 Advancing Big Data Science in Genomics Research program, which is jointly funded by the
334 Natural Sciences and Engineering Research Council (NSERC) of Canada, the Canadian
335 Institutes of Health Research (CIHR), Genome Canada and the Canada Foundation for
336 Innovation (CFI). P.C.B. was supported by a Terry Fox Research Institute New Investigator Award
337 and a CIHR New Investigator Award. This work was supported by an NSERC Discovery grant
338 and by Canadian Institutes of Health Research, grant #SVB-145586, to PCB. The results
339 described here are in part based upon data generated by the ICGC/TCGA Pan-Cancer Analysis
340 of Whole Genomes Network: <https://dcc.icgc.org/>

341 Author Contributions

342 CHL and PCB initiated the project. CHL, SDP, RXS, FY and NS analyzed data. PCB supervised
343 research. CHL and PCB wrote the first draft of the manuscript, which all authors edited and
344 approved. The PCAWG network provided variant calls and insightful commentary.

345 References (Main Text)

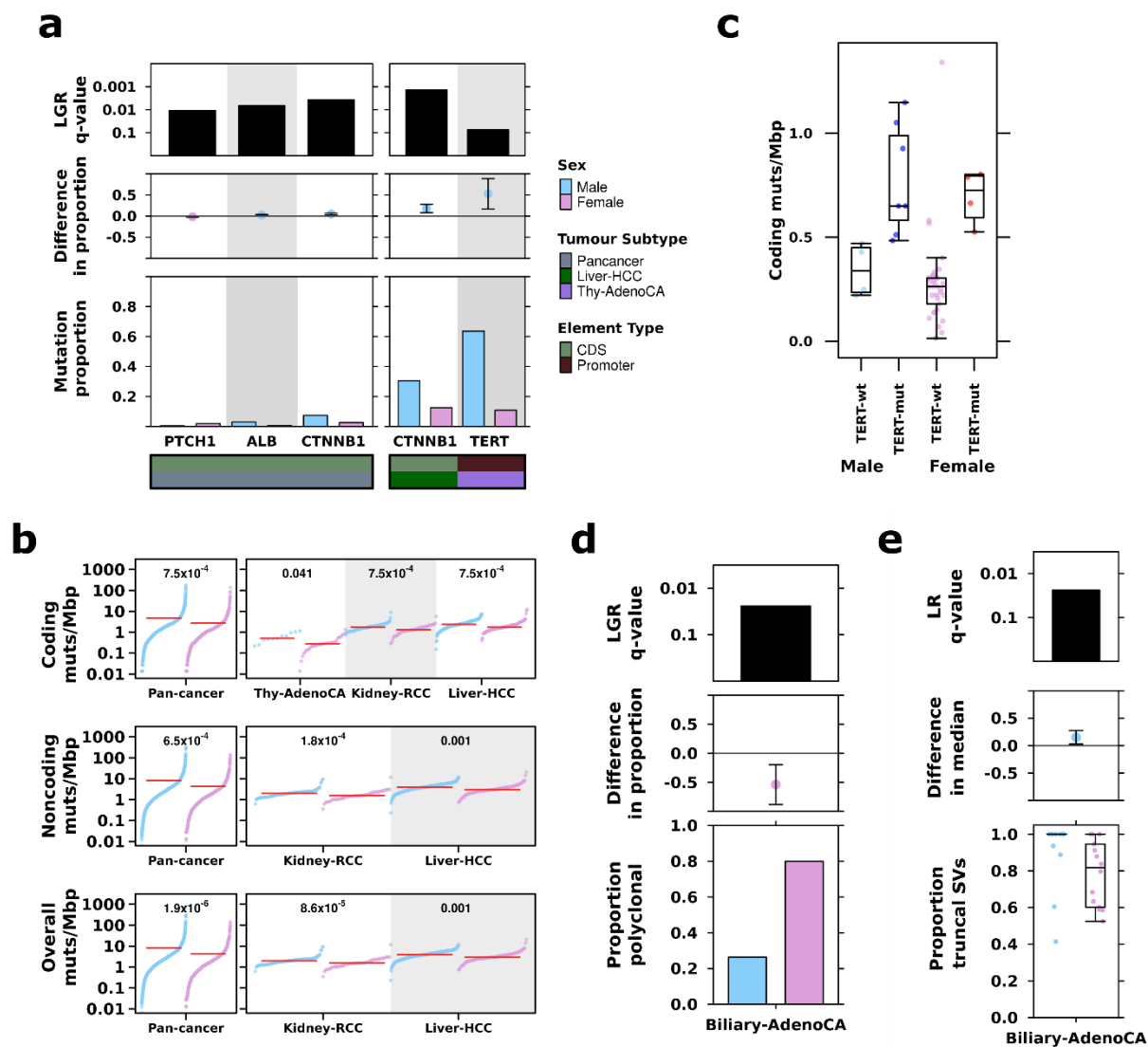
- 346 1. Cook, M. B. *et al.* Sex disparities in cancer incidence by period and age. *Cancer Epidemiol.*
347 *Biomark. Prev. Publ. Am. Assoc. Cancer Res. Cosponsored Am. Soc. Prev. Oncol.* **18**, 1174–
348 1182 (2009).
- 349 2. Edgren, G., Liang, L., Adami, H.-O. & Chang, E. T. Enigmatic sex disparities in cancer
350 incidence. *Eur. J. Epidemiol.* **27**, 187–196 (2012).
- 351 3. Cook, M. B., McGlynn, K. A., Devesa, S. S., Freedman, N. D. & Anderson, W. F. Sex
352 Disparities in Cancer Mortality and Survival. *Cancer Epidemiol. Biomarkers Prev.* **20**, 1629–
353 1637 (2011).
- 354 4. Wichmann, M. W., Muller, C., Hornung, H. M., Lau-Werner, U. & Schildberg, F. W. Gender
355 differences in long-term survival of patients with colorectal cancer. *Br. J. Surg.* **88**, 1092–1098
356 (2001).
- 357 5. Elsaleh, H. *et al.* Association of tumour site and sex with survival benefit from adjuvant
358 chemotherapy in colorectal cancer. *Lancet Lond. Engl.* **355**, 1745–1750 (2000).
- 359 6. OuYang, P.-Y. *et al.* The significant survival advantage of female sex in nasopharyngeal
360 carcinoma: a propensity-matched analysis. *Br. J. Cancer* **112**, 1554–1561 (2015).
- 361 7. Gupta, S., Artomov, M., Goggins, W., Daly, M. & Tsao, H. Gender Disparity and Mutation
362 Burden in Metastatic Melanoma. *J. Natl. Cancer Inst.* **107**, djv221 (2015).
- 363 8. Sun, T., Plutynski, A., Ward, S. & Rubin, J. B. An integrative view on sex differences in brain
364 tumors. *Cellular and Molecular Life Sciences* **72**, 3323–3342 (2015).
- 365 9. Dunford, A. *et al.* Tumor-suppressor genes that escape from X-inactivation contribute to
366 cancer sex bias. *Nat. Genet.* (2016). doi:10.1038/ng.3726
- 367 10. Yuan, Y. *et al.* Comprehensive Characterization of Molecular Differences in Cancer between
368 Male and Female Patients. *Cancer Cell* **29**, 711–722 (2016).

- 369 11. Warrington, N. M. *et al.* The Cyclic AMP Pathway Is a Sex-Specific Modifier of Glioma Risk
370 in Type I Neurofibromatosis Patients. *Cancer Res.* **75**, 16–21 (2015).
- 371 12. Li, C. H., Haider, S., Shiah, Y.-J., Thai, K. & Boutros, P. C. Sex Differences in Cancer Driver
372 Genes and Biomarkers. *Cancer Res* **78**, 5527 (2018).
- 373 13. Lopes-Ramos, C. M. *et al.* Gene Regulatory Network Analysis Identifies Sex-Linked
374 Differences in Colon Cancer Drug Metabolism. *Cancer Research* **78**, 5538–5547 (2018).
- 375 14. Campbell, P. J., Getz, G., Stuart, J. M., Korb, J. O. & Stein, L. D. Pan-cancer analysis of
376 whole genomes. *bioRxiv* (2017). doi:10.1101/162784
- 377 15. Dunford, A. *et al.* Tumor-suppressor genes that escape from X-inactivation contribute to
378 cancer sex bias. *Nature Genetics* (2016). doi:10.1038/ng.3726
- 379 16. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *bioRxiv* (2017).
380 doi:10.1101/190330
- 381 17. Liu, R. & Xing, M. TERT promoter mutations in thyroid cancer. *Endocrine-Related Cancer*
382 **ERC-15-0533** (2016). doi:10.1530/ERC-15-0533
- 383 18. Qin, Y. *et al.* The inhibitory effects of deleted in liver cancer 1 gene on gallbladder cancer
384 growth through induction of cell cycle arrest and apoptosis: Deleted in liver cancer 1 in
385 gallbladder cancer. *Journal of Gastroenterology and Hepatology* **29**, 964–972 (2014).
- 386 19. Nik-Zainal, S. *et al.* Mutational Processes Molding the Genomes of 21 Breast Cancers. *Cell*
387 **149**, 979–993 (2012).
- 388 20. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,
389 415–421 (2013).
- 390 21. Alexandrov, L. *et al.* The Repertoire of Mutational Signatures in Human Cancer. (2018).
391 doi:10.1101/322859
- 392 22. Cho, A., Chang, Y., Ahn, J., Shin, H. & Ryu, S. Cigarette smoking and thyroid cancer risk: a
393 cohort study. *British Journal of Cancer* **119**, 638–645 (2018).

- 394 23. Yeo, I.-K. & Johnson, R. A. A New Family of Power Transformations to Improve Normality or
395 Symmetry. *Biometrika* 87, 954–959 (2000).
- 396 24. P'ng, C. et al. BPG: Seamless, automated and interactive visualization of scientific data. *BMC*
397 *Bioinformatics* 20, (2019).
- 398
- 399

400 Figures

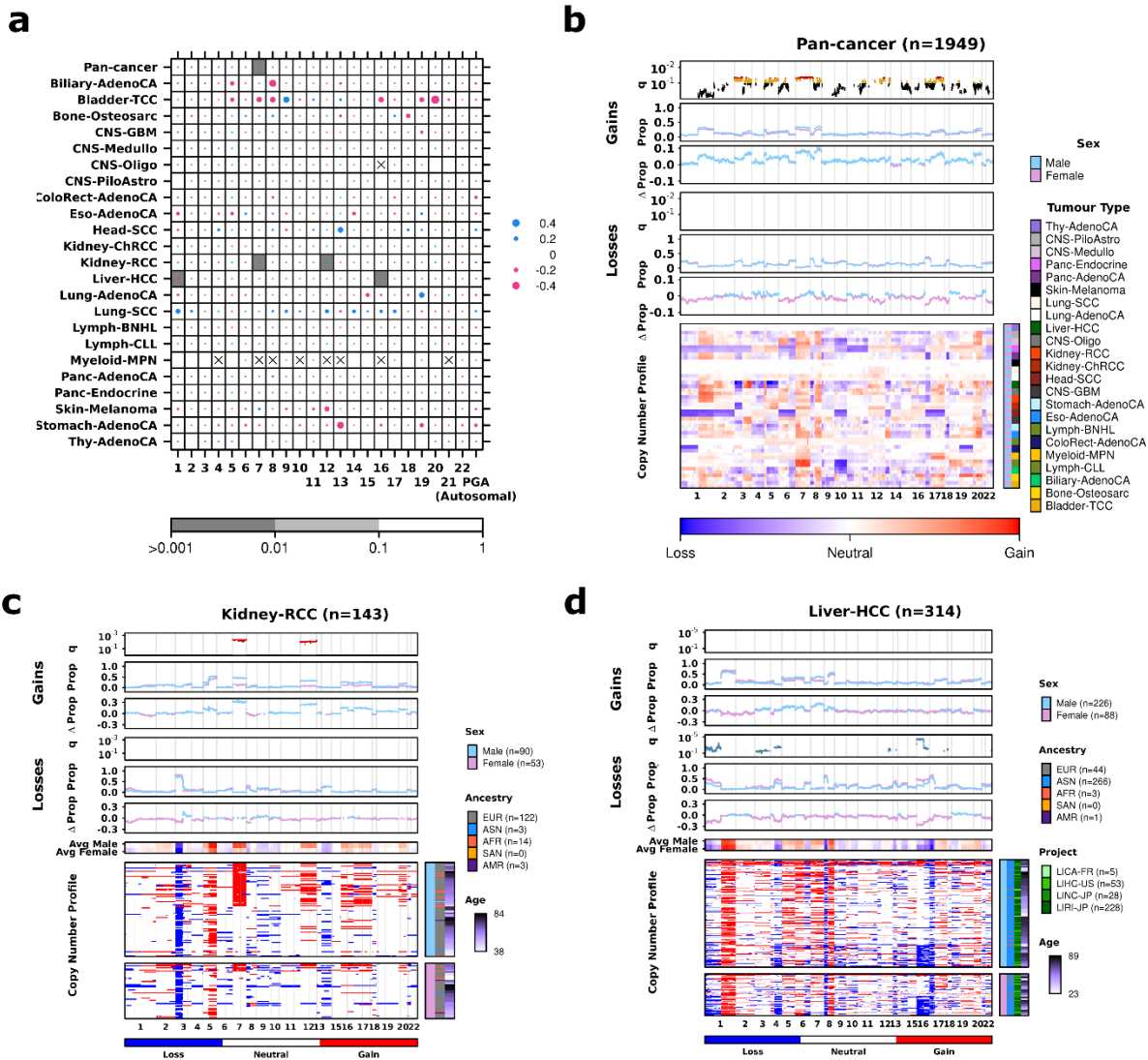
401 Figures 1-4



402

403 **Figure 1 | Sex biases in mutation frequency of driver genes, mutation prevalence and**
 404 **tumour evolution. (a)** From top to bottom, each plot shows the logistic regression q-value for the
 405 sex effect; difference in proportion of mutated samples between the sexes, where blue denotes
 406 male-bias and pink denotes female-bias; and mutation proportion for each gene. Bottom covariate
 407 bars indicate mutation context and tumour subtype of interest. **(b)** The burden of somatic SNVs
 408 for coding, non-coding and overall mutation load. Linear regression q-values are shown. **(c)**
 409 Coding mutation load for thyroid adenocarcinoma samples compared by sex and presence or
 410 absence of TERT promoter mutations. **(d)** The proportion of polyclonal samples and **(e)** the
 411 proportion of truncal structural variants in biliary cancer.
 412

413

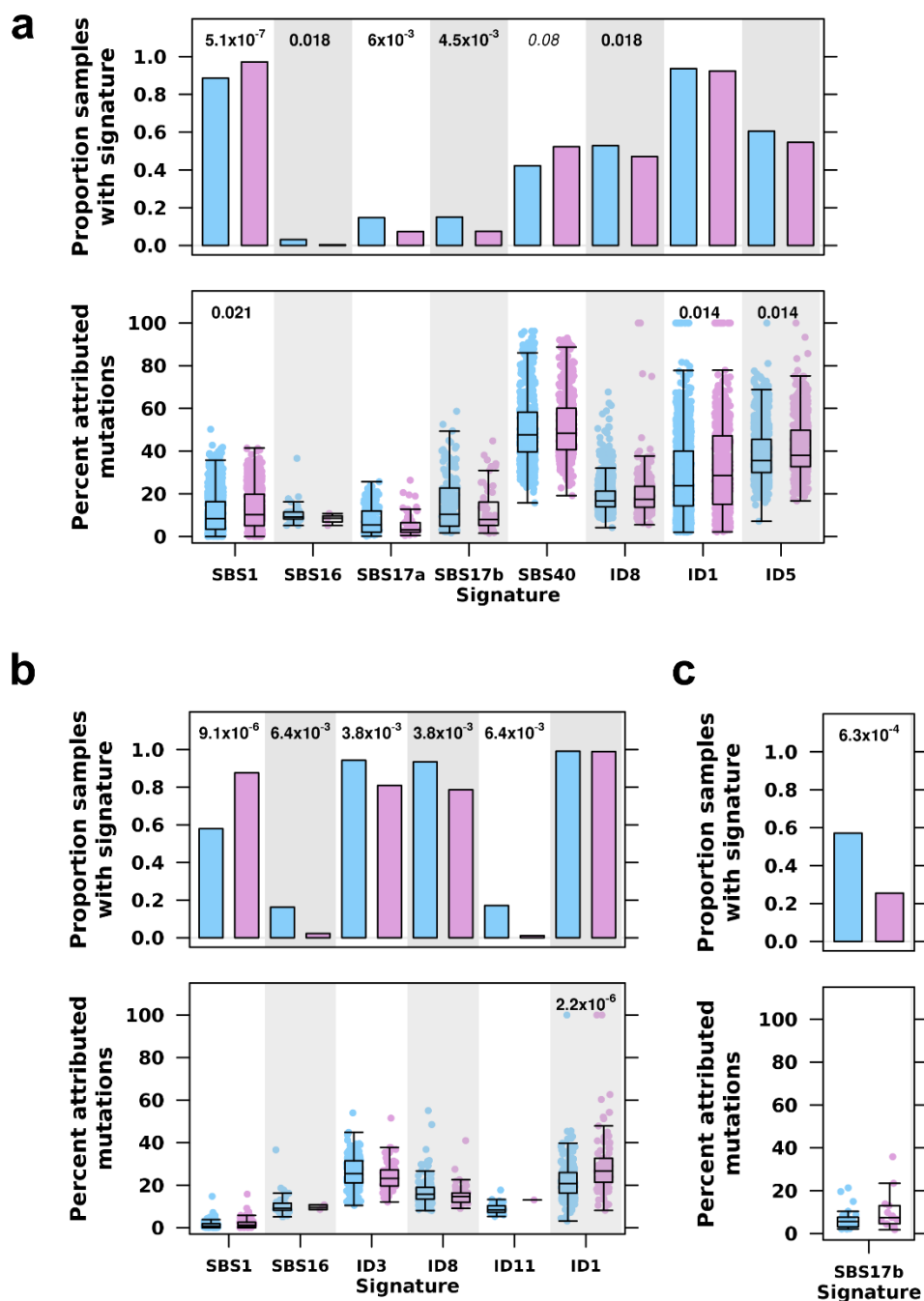


414

415 **Figure 2 | Sex-biases in percent chromosome altered are reflected in gene-specific events.**

416 **(a)** Dotmap showing association between sex and percent genome or chromosome altered,
 417 where dot size shows difference in median percent genome/chromosome altered between the
 418 sexes, and background shading shows q-values from linear regression. Sex differences in CNAs
 419 for **(b)** pan-cancer, **(c)** kidney renal cell cancer and **(d)** hepatocellular cancer. Each plot shows,
 420 from top to bottom: the q-value showing significance of sex from multivariate linear modeling with
 421 yellow (green) points corresponding to $0.1 < q < 0.05$ and deep blue (red) points corresponding
 422 to $q < 0.05$; the proportion of samples with aberration; the difference in proportion between male
 423 and female groups for copy number gain events; the same repeated for copy number loss events;
 424 and the copy number aberration (CNA) profile heatmap. The columns represent genes ordered
 425 by chromosome. Light blue and pink points represent data for male- and female- derived samples
 426 respectively.

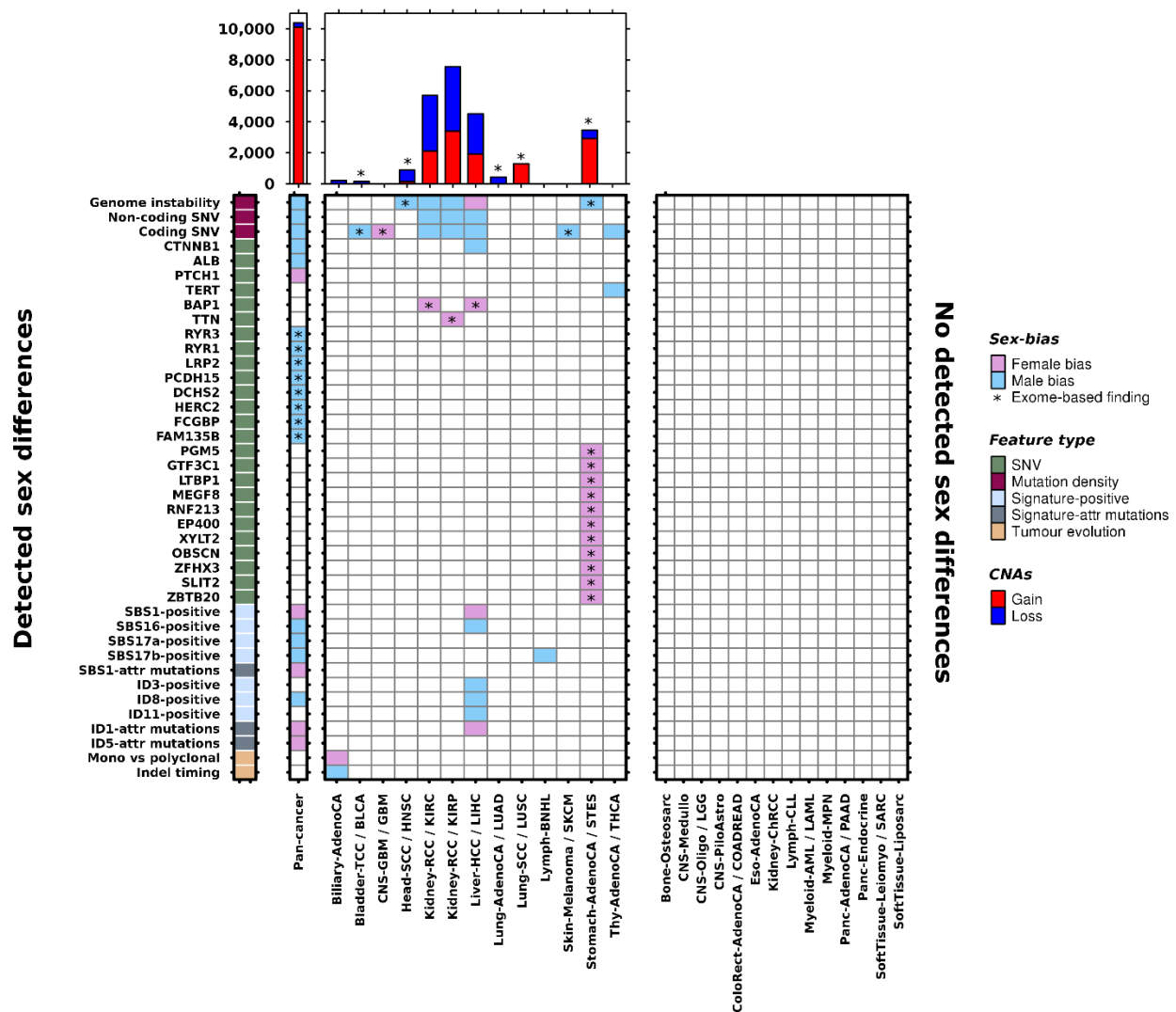
427



428

429 **Figure 3 | Sex differences in trinucleotide signatures related to mutational processes.**

430 Comparisons between proportion of signature positive samples shown in barcharts and proportion
 431 of attributed mutations shown in boxplots for **(a)** pan-cancer comparisons, **(b)** liver hepatocellular
 432 cancer, and **(c)** B-cell non-Hodgkin lymphoma. FDR-adjusted q-values for logistic regression
 433 (barplots) and linear regression (boxplots) shown only for significant comparisons. Blue shows
 434 male- and pink shows female-derived samples.



435

436 **Figure 4 | The landscape of sex differences in cancer genomics.** Heatmap shows genomic
 437 features found to be sex-biased in pan-cancer analysis or in specific tumour subtypes. Results
 438 from both PCAWG and TCGA analyses are shown. Direction of sex-bias is shown in coloration
 439 denoting which sex has higher or more frequent aberration of the genomic feature. Top barplot
 440 shows union of genes found to be involved in sex-biased CNAs. Starred indicate findings
 441 exclusively from exome sequencing data (n=7,131) and unstarred indicate findings from PCAWG
 442 data (n=1,983).

443

444

445 Online Methods

446 Data acquisition & Processing

447 Data was downloaded from the PCAWG consortium through Synapse. All data pre-processing
448 was performed by the consortium as described¹⁴. Additional data-specific details are described
449 below.

450 General Statistical Framework

451 We followed a statistical approach as previously described in our previous work (Li et al). Briefly,
452 for each genomic feature of interest, we used univariate tests first followed by false discovery rate
453 (FDR) adjustment to identify putative sex-biases of interest ($q < 0.1$). Here, we use non-parametric
454 univariate tests to minimize assumptions on the data. For putative sex-biases, we then follow up
455 the univariate analysis with multivariate modeling to account for potential confounders using
456 bespoke models for each tumour subtype. Model variables for each tumour context are described
457 in **Supplementary Table 1** and were included based on availability of data (<15% missing),
458 sufficient variability (at least two levels) and collinearity. Discrete data was modeled using logistic
459 regression. Continuous data was first transformed using the Box-Cox family and modeled using
460 linear regression. The Box-Cox family of transformations is a formalized method to select a power
461 transformation to better approximate a normal-like distribution and stabilize variance. We used
462 the Yeo-Johnson extension to the Box-Cox transformation that allows for zeros and negative
463 values²³:

$$464 \quad y_i^\lambda = \begin{cases} \frac{(y_i + 1)^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, y \geq 0 \\ \log(y_i + 1), & \text{if } \lambda = 0, y \geq 0 \\ -\frac{(-y_i + 1)^{2-\lambda} - 1}{2 - \lambda}, & \text{if } \lambda \neq 2, y < 0 \\ -\log(-y_i + 1), & \text{if } \lambda = 2, y < 0 \end{cases}$$

465 FDR adjustment was performed for p-values for the sex variable significance estimate and an
466 FDR threshold of 10% was used to determine statistical significance. More detail is provided for
467 each analysis below.

468 Driver Event Analysis

469 We focused on driver events described by the PCAWG consortium¹⁶. Driver mutation data was
470 binarized to indicate presence or absence of the driver event in each patient. Proportions of
471 mutated genes were compared between the sexes using proportions tests for univariate analysis.
472 A q-value threshold of 0.1 was used to select genes for further multivariate analysis using binary
473 logistic regression. FDR correction was again applied and genes with significant pan-cancer sex
474 terms were extracted from the models (q-value < 0.1). Driver event analysis was performed
475 separately for pan-cancer analysis and for each tumour subtype.

476 **Clonal structure and mutation timing analysis**

477 Subclonal structure and mutation timing calls were downloaded from Synapse
478 (<https://www.synapse.org/#!/Synapse:syn8532425>). Subclonal structure data was binarized from
479 number of subclonal clusters per sample to monoclonal (one cluster) or polyclonal (more than
480 one cluster). The proportion of polyclonal samples was calculated per sex and compared using
481 proportion tests for both pan-cancer and tumour subtype analysis. The univariate p-values were
482 FDR adjusted across all tumour subtypes to identify putatively sex-biased clonal structure. These
483 cases were further scrutinized using logistic regression. A multivariate q-value threshold of 0.1
484 was used to determine statistically significant sex-biased clonal structure.

485 Mutation timing data classified SNVs, indels and SVs into clonal (truncal) or subclonal groups.
486 The proportion of truncal variants was calculated for each mutation type ($\frac{\text{Number truncal SNVs}}{\text{total SNVs}}$, etc)
487 to obtain proportions of truncal SNVs, indels and SVs for each sample. These proportions were
488 compared between the sexes using two-sided Mann-Whitney *U*-tests and univariate p-values
489 were FDR adjusted to identify putatively sex-biased mutation timing. Linear regression was used
490 to adjust for confounding factors and a multivariate q-value threshold of 0.1 was used to determine
491 statistically significant sex-biased mutation timing. The mutation timing analysis was performed
492 separately for SNVs, indels and SVs.

493 **Mutation Load analysis**

494 Consensus SNV calls were downloaded from Synapse
495 (<https://www.synapse.org/#!/Synapse:syn7118450>). Overall mutation prevalence per patient was
496 calculated as the sum of SNVs across all genes on the autosomes and scaled to mutations/Mbp.
497 Coding mutation prevalence only considers the coding regions of the genome, and noncoding
498 prevalence only considers the noncoding regions. Mutation load was compared between the
499 sexes using Mann-Whitney *U*-tests for both pan-cancer and tumour-type specific analysis.
500 Comparisons with u-test q-values meeting an FDR threshold of 10% were further analyzed using
501 linear regression to adjust for tumour subtype-specific variables. Mutation load analysis was
502 performed separately for each mutation context, with pan-cancer and tumour subtype p-values
503 adjusted together

504 **Chromosome and Genome Instability analysis**

505 Consensus copy number data was obtained from Synapse
506 (<https://www.synapse.org/#!/Synapse:syn8042880>). Ploidy-adjusted calls were used to identify
507 segments with copy number gains and losses. The number of bases in copy number gained or
508 lost segments were summed per chromosome and divided by chromosome size to obtain percent
509 chromosome gained and lost, respectively. All segments affected by a copy number aberration
510 were also summed and treated in the same way to calculate percent chromosome altered.
511 Percent copy number gained, lost, and altered were also calculated over the autosomes. These
512 metrics were compared in pan-cancer and tumour-subtype analysis using u-tests to identify
513 putatively sex-biased chromosome and genome instability, and putatively sex-biased events were
514 further analysed using linear regression modeling. Genome instability analysis was performed
515 separately for each tumour subtype with FDR adjustment performed over percent copy gained,
516 loss and altered comparisons together.

517 **Genome-spanning CNA analysis**

518 Consensus copy number data was processed to gain/neutral/loss calls per gene. The number of
519 loss, neutral and gain calls were summed per sex, and assessed using univariate and multivariate
520 techniques. For univariate analysis, proportional differences between the sexes for gains and
521 losses were tested for each gene using proportions tests. After identifying candidate pan-cancer
522 univariately significant genes, multivariate logistic regression was used to adjust ternary CNA data
523 for sex, age, ancestry and tumour-type. The genome-spanning analysis was performed
524 separately for losses and gains for each tumour subtype.

525 **Mutational Signatures analysis**

526 The number of mutations attributed to each SBS, DBS and ID signature per sample was
527 downloaded from Synapse (<https://www.synapse.org/#!/Synapse:syn8366024>). For each
528 signature, we compared the proportion of samples with any mutations attributed to the signatures
529 (“signature-positive”) using χ^2 -square tests to identify univariately significant sex-biases.
530 Signatures with putative sex-biases were further analysed using logistic regression.

531 We also compared the proportions of mutations attributed to each signature. The numbers of
532 mutations per signature were divided by total number of mutations for each sample to obtain the
533 proportion of mutations attributed to the signature. Mann-Whitney *U*-tests were used to compare
534 these proportions. Putative sex-biased signatures were further analysed using linear regression
535 after Box-cox adjustment.

536 Signatures that were not detected in a tumour subtype was omitted from analysis for that tumour
537 subtype. Statistical analyses were performed for each set tumour subtype, but combining all SBS,
538 DBS and ID signatures.

539 **Statistical Analysis & Data Visualization**

540 All statistical analyses and data visualization were performed in the R statistical environment
541 (v3.4.3) using the BPG²⁴ (v5.9.8), car (v3.0-2) and mlogit (v0.2-4), packages.