

African-specific prostate cancer molecular taxonomy

Vanessa Hayes (✉ vanessa.hayes@sydney.edu.au)

University of Sydney <https://orcid.org/0000-0002-4524-7280>

Weerachai Jaratlerdsiri

University of Sydney

Jue Jiang

Garvan Institute of Medical Research <https://orcid.org/0000-0003-0920-8310>

Tingting Gong

University of Sydney

Sean Patrick

University of Pretoria

Cali Willet

University of Sydney

Tracy Chew

University of Sydney

Ruth Lyons

Garvan Institute of Medical Research

Anne-Maree Haynes

Garvan Institute of Medical Research

Gabriela Pasqualim

Universidade Federal do Rio Grande do Sul

Melanie Louw

National Health Laboratory Services

James Kench

University of Sydney

Raymond Campbell

Kalafong Academic Hospital

Eva Chan

New South Wales Health Pathology <https://orcid.org/0000-0002-6104-3763>

David Wedge

University of Manchester <https://orcid.org/0000-0002-7572-3196>

Rosemarie Sadsad

University of Sydney

Ilma Brum

Universidade Federal do Rio Grande do Sul

Shingai Mutambirwa

Sefako Makgatho Health Science University

Phillip Stricker

St Vincnet's Hospital

Riana Bornman

University of Pretoria <https://orcid.org/0000-0003-3975-2333>

Lisa Horvath

Chris O'Brien Lifehouse

Biological Sciences - Article

Keywords: Prostate cancer, tumour genome profiling, Global Mutational Subtypes

Posted Date: December 1st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1122619/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature on August 31st, 2022. See the published version at <https://doi.org/10.1038/s41586-022-05154-6>.

1 **African-specific prostate cancer molecular taxonomy**

2

3 Weerachai Jaratlerdsiri^{1,2}, Jue Jiang^{1,2}, Tingting Gong^{1,2}, Sean M. Patrick³, Cali Willet⁴,
4 Tracy Chew⁴, Ruth J. Lyons², Anne-Maree Haynes⁵, Gabriela Pasqualim^{6,7}, Melanie
5 Louw⁸, James G. Kench⁹, Raymond Campbell¹⁰, Lisa G. Horvath^{5,11}, Eva K.F. Chan²,
6 David C. Wedge¹², Rosemarie Sadsad⁴, Ilma Simoni Brum⁶, Shingai B.A.
7 Mutambirwa¹³, Phillip D. Stricker^{5,14}, M.S. Riana Bornman³, Vanessa M. Hayes^{1,2,3,15*}

8

9 ¹Ancestry and Health Genomics Laboratory, Charles Perkins Centre, School of Medical
10 Sciences, Faculty of Medicine and Health, University of Sydney, Camperdown, NSW,
11 Australia; ²Human Comparative and Prostate Cancer Genomics Laboratory, Garvan Institute of
12 Medical Research, Darlinghurst, NSW, Australia; ³School of Health Systems & Public Health,
13 University of Pretoria, South Africa; ⁴Sydney Informatics Hub, University of Sydney,
14 Darlington, NSW, Australia; ⁵Genomics and Epigenetics Theme, Garvan Institute of Medical
15 Research, Darlinghurst, NSW, Australia; ⁶Endocrine and Tumor Molecular Biology Laboratory
16 (LABIMET), Instituto de Ciências Básicas da Saúde, Universidade Federal do Rio Grande do
17 Sul, Brazil; ⁷Laboratory of Genetics, Instituto de Ciências Biológicas, Universidade Federal do
18 Rio Grande, Brazil; ⁸National Health Laboratory Services, Johannesburg, South Africa;
19 ⁹Department of Tissue Pathology and Diagnostic Oncology, Royal Prince Alfred Hospital and
20 Central Clinical School, University of Sydney, Sydney, NSW, Australia; ¹⁰Kalafong Academic
21 Hospital, Pretoria, South Africa; ¹¹Medical Oncology, Chris O'Brien Lifecare, Royal Prince
22 Alfred Hospital and Faculty of Medicine and Health, University of Sydney Camperdown,
23 NSW, Australia; ¹²Division of Cancer Sciences, University of Manchester, United Kingdom;
24 ¹³Department of Urology, Sefako Makgatho Health Science University, Dr George Mukhari
25 Academic Hospital, Medunsa, South Africa; ¹⁴Department of Urology, St. Vincent's Hospital,

26 Darlinghurst, NSW, Australia; ¹⁵Faculty of Health Sciences, University of Limpopo, Turfloop
27 Campus, South Africa.

28 *e-mail: vanessa.hayes@sydney.edu.au

29

30 **Abstract**

31 Prostate cancer is characterised by significant global disparity; mortality rates in Sub-
32 Saharan Africa are double to quadruple those in Eurasia¹. Hypothesising unknown
33 interplay between genetic and non-genetic factors, tumour genome profiling
34 envisages contributing mutational processes^{2,3}. Through whole-genome sequencing of
35 treatment-naïve prostate cancer from 183 ethnically/globally distinct patients (African
36 *versus* European), we generate the largest cancer genomics resource for Sub-Saharan
37 Africa. Identifying ~2 million somatic variants, Africans carried the greatest burden.
38 We describe a new molecular taxonomy using all mutational types and ethno-
39 geographic identifiers, including Asian. Defined as Global Mutational Subtypes
40 (GMS) A–D, although Africans presented within all subtypes, we found GMS-B to be
41 ‘African-specific’ and GMS-D ‘African-predominant’, including Admixed and
42 European Africans. Conversely, Europeans from Australia, Africa and Brazil
43 predominated within ‘mutationally-quiet’ and ethnically/globally ‘universal’ GMS-A,
44 while European Australians shared a higher mutational burden with Africans in GMS-
45 C. GMS predicts clinical outcomes; reconstructing cancer timelines suggests four
46 evolutionary trajectories with different mutation rates (GMS-A, low 0.968/year *versus*
47 D, highest 1.315/year). Our data suggest both common genetic factors across extant
48 populations and regional environmental factors contributing to carcinogenesis,
49 analogous to gene-environment interaction defined here as a different effect of an
50 environmental surrounding in persons with different ancestries or vice versa. We

- 51 anticipate GMS acting as a proxy to intrinsic and extrinsic mutational processes in
- 52 cancers, promoting global inclusion in landmark studies.

53 **Main**

54 Prostate cancer is a common heterogeneous disease, responsible annually for more
55 than 1,400,000 new diagnoses and 375,000 male-associated deaths worldwide¹.
56 Characterised by a highly variable natural history and diverse clinical behaviours⁴, it
57 is not surprising that genome profiling has revealed extensive intra- and inter-tumour
58 heterogeneity and complexity^{5,6}. The identification of oncogenic subtypes⁷ and
59 actionable drug targets⁸ are moving prostate cancer management a step closer to the
60 promise of precision medicine^{7,9-13}. While high-income European ancestral countries
61 are well along the road to incorporating cancer genomics in all aspects of cancer
62 care¹⁴, the rest of the world lags behind, with a notable absence in Sub-Saharan
63 Africa¹⁵. Prostate cancer is no different, with a single large-scale study out of China¹²;
64 in 2018, we provided the first snapshot for Sub-Saharan Africa, reporting an elevated
65 mutational density in a mere six cases¹⁶. With mortality rates over double high-
66 income countries and quadrupled for greater Asia, Sub-Saharan Africa prostate cancer
67 is the top-ranked male-associated cancer both by diagnosis and deaths, including
68 southern Africa with age-standardised rates of 65.9 and 22 per 100,000, respectively¹.
69 Through the Southern African Prostate Cancer Study (SAPCS), we report a 2.1-fold
70 increase in aggressive disease compared to African Americans¹⁷.

71 Here we describe, to our knowledge, the largest cancer and prostate cancer genomics
72 data for Sub-Saharan Africa, including 123 South African men. Controlling for study
73 artefacts, an additional 60 non-Africans were passed simultaneously through the same
74 high-depth whole-genome sequencing (WGS), mutation-calling and analytical
75 framework. Focusing on treatment-naïve aggressive tumours (mostly Grades 4-5,
76 Extended Data Fig. 1a) and patient-matched blood achieving coverages of

77 88.69±14.78 and 44.34±8.11, respectively (median±s.d., Supplementary Table 1), we
78 uniformly generated, called and assessed about 2 million somatic variants. We show a
79 greater number of acquired genetic alterations within Africans, while identifying both
80 globally relevant and African-specific genomic subtypes. Through combining our
81 somatic variant dataset with that published for European-ancestral^{7,8,18,19} and Chinese¹²
82 prostate cancer genomes, we reveal a novel prostate cancer taxonomy with different
83 clinical outcomes. The inclusion of 2,658 cancer genomes from the ICGC/TCGA
84 Pan-Cancer Analysis of Whole Genomes (PCAWG)¹⁴ led to expanding our global
85 mutational subtyping between cancer types. Using known clock-like mutational
86 processes in each subtype, we infer mutation timing of oncogenic drivers in broad
87 periods of tumour evolution and calculate mutation rates for each subtype that had a
88 distinctive tumour evolution pattern. Combined, these analyses allow us to
89 demonstrate how global inclusion in cancer genomics can unravel unseen
90 heterogeneity in prostate cancer in terms of its genomic and clinical behaviours.

91 **Genetic ancestry**

92 Genetic ancestries were estimated for the 183 patient donors using a joint dataset in a
93 unified analysis aggregated from a collection of geographically matched African
94 (n=64) and European (n=4) deep-coverage reference genomes^{20,21}. Ancestries were
95 assigned using 7,472,833 markers as: African (n=113), with greater than 98%
96 contribution; European (n=61), allowing for up to 10% Asian contribution (with a
97 single outlier of 26%); and African-European Admixed (n=9), with as little as 4%
98 African or European contribution (Extended Data Fig. 1b).

99 **Total somatic mutations**

100 In 183 prostate tumours, we identified 1,067,885 single nucleotide variants (SNVs),
101 11,259 dinucleotides, 307,263 small insertions/deletions (indels <50 bp), 419,920
102 copy number alterations (CNAs) and 22,919 structural variants (SVs), with each
103 mutational type elevated in African derived tumours (Fig. 1a). A median of
104 37.54%±5.51 of SNVs were C-to-T mutations, and the transition and transversion
105 ratio was 1.282 cohort-wise. African derived tumours harboured a higher rate of small
106 mutations (SNVs and indels), with a median of 1.197 mutations/Mb (0.031-170.445),
107 compared to those of Europeans (1.061 mutations/Mb, *P*-value = 0.013, two-sample t-
108 test). Percent genome alteration (PGA) was similarly greater in Africans (0.073 *versus*
109 0.028, *P*-value = 0.021). Correlation tests of ethnicity and total somatic mutations also
110 supported the findings (FDR=0.009 and 0.032 for SNVs and PGA, respectively,
111 Extended Data Fig. 1d). The top six highest estimates of SV breakpoints per sample
112 were observed among African patients (928-2,284 breakpoints). Intrachromosomal
113 SV breakpoints were 52-55% positive for chromothripsis among Africans and
114 Europeans (median, 3 and 2 high-confidence events, respectively). Chromoplexy was
115 more frequent in Europeans than in Africans (38% *versus* 33%, *P*-value=0.536), with
116 the number of interchromosomal chains more likely to be elevated in Africans than
117 Europeans (1-6 *versus* 1-2, *P*-value=0.748). Moreover, the magnitude of all types of
118 mutations was strongly correlated to one another (Fig. 1b). Thus, the more mutations
119 a prostate tumour has of any given type, the more mutations it is likely to have of all
120 types.

121 **Candidate oncogenic drivers**

122 Prostate cancer is known to have a long tail of oncogenic drivers¹⁹ across the
123 spectrum of different mutational types⁸ (Extended Data Fig. 2). Protein-coding

124 mutations, including probably and possibly damaging, were significantly greater in
125 Africans (PolyPhen-2, 14 *versus* 11 mutations in Europeans, P -value=0.022, two-
126 sample t-test). We identified 482 coding and 167 noncoding drivers defined by the
127 PCAWG Consortium²² (Extended Data Fig. 3a). A median of 2 ± 22.5 coding drivers
128 was observed in this study (Supplementary Table 2), with 1 ± 5.4 appearing to be
129 prostate cancer-specific^{7,8,18,19}. The coding driver genes significantly mutated among
130 183 patients were *FOXA1*, *PTEN*, *SPOP* and *TP53* (10-25 patients, FDR=1.34e-21–
131 9.44e-05), while noncoding driver elements were the *FOXA1* 3'-UTR, *SNORD3B-2*
132 small RNA and a regulatory miRNA promoter at chromosome 22:38,381,983
133 (FDR=9.12e-13, 6.16e-09 and 0.070, respectively). Recurrent CNAs of all the
134 patients included 137 gains and 129 losses (*GISTIC2*, FDR <0.10, Supplementary
135 Table 3) with some spanning driver genes (Extended Data Fig. 3b), such as *DNAH2*
136 (FDR=2.18e-07), *FAM66C* (1.30e-09), *FOXP1* (0.005), *FXR2* (2.18e-07), *PTEN*
137 (9.61e-13), *SHBG* (2.18e-07), and *TP53* (2.18e-07).

138 In addition, a fraction of somatic SVs (2 breakpoints each; 1,328 breakpoints in total)
139 overlapped with 156 driver genes reported as altered by significantly recurrent
140 breakpoints in the PCAWG study²², while using a generalised linear model with
141 adjustable background covariates we identified an additional 100 genes to be
142 significantly impacted by SV breakpoints (FDR=1.3e-43–0.097, Extended Data Fig.
143 3c, Supplementary Table 4). For over 20% of tumours, SV breakpoints coexisted with
144 other mutational types within *DNAH2*, *ERG*, *FAM66C*, *FXR2*, *PTEN*, *SHBG*, and
145 *TP53*. Using optical genome mapping (OGM), an alternative non-sequencing method
146 to interrogate for chromosomal abnormalities²³, we validated recurrent breakpoints in
147 novel HLA regions (*DQA1* and *DQB1* genes), identifying translocations between the

148 3-Mb HLA complex at chromosome 6 and its corresponding HLA alternate contigs
149 (Extended Data Fig. 3d).

150 **Integrative clustering analysis of molecular subtypes**

151 Molecular subtyping of tumours is a standard approach in cancer genomics to stratify
152 patients into different degrees of somatic alterations in a homogeneous population,
153 with an implication for clinical use⁹⁻¹². Identifying five of the seven TCGA oncogenic
154 driver-defined subtypes in our study⁷, European patients were 25% more likely than
155 African patients to be classified (Supplementary Table 5, Extended Data Fig. 4a-d).
156 While *TMPRSS2-ERG* fusions (predominantly 3-Mb deletions) and *FOXA1* coding
157 mutations (forkhead domain) occurred at higher frequencies in our European over
158 African patients, 37.7% and 8.2% *versus* 13.3% and 5.3%, respectively (OR=0.255,
159 *P*-value=0.0004 and OR=0.854, *P*-value=0.771), *SPOP* coding mutations (MATH
160 and BTB domains) were more common in the African (8.8%) *versus* European
161 patients (6.6%, OR=1.688, *P*-value=0.426).

162 For further molecular classification, we performed iCluster analysis on all mutational
163 types (small mutations, copy number and SVs) identifying four subtypes, A to D
164 (Supplementary Table 6, Fig. 2a, b). We found Subtype A to be mutationally quiet
165 (1.01 mutations/Mb, 0.50 breakpoints/10Mb, 2% PGA); conversely Subtype D
166 showed the greatest mutational density (1.91 mutations/Mb, 1.08 breakpoints/10Mb,
167 31% PGA) with a mixture of copy number (CN) gains and losses, while Subtypes B
168 and C were marked by substantial CN gains or losses, respectively (Fig. 3b). The
169 quiet subtype seems to be common in prostate cancer studies^{7,9,24}, while the number
170 of pan-cancer consensus drivers²² increased from Subtype A (median, 2 drivers) to B
171 (3), C (3) and D (4).

172 Using all mutational types in the analysis, 124 genes were significantly mutated
173 across the four subtypes (FDR=3.742e-13–0.067; Fig. 3a), occurring in 31 to 183
174 patients (frequency, 0.17-1). Among them, 100 genes were reported as oncogenic
175 drivers in the PCAWG²² and *FOXA1* and *SPOP* genes acting as the TCGA subtypes
176 were also replicated in this analysis, while the 24 novel recurrently mutated genes
177 were predominantly impacted by SV breakpoints and CNAs. The median number of
178 mutated genes ranged from 28 (range 3-105) for Subtype A to 82, 98 and 93 for
179 Subtypes B, C and D, respectively (42-109, 72-112, 49-107). While different
180 mutational types tended to co-occur within genes and/or patients (Supplementary
181 Table 7), small mutations (coding and noncoding) were noticeably observed in the
182 quiet subtype, supporting acquisition early in tumorigenesis²⁵. Our preferentially
183 mutated genes within tumour subtypes resemble the long tail of prostate cancer
184 drivers¹⁹, with some highly impacting many tumours, but most only impacting a few
185 tumours.

186 The 124 preferentially mutated genes within our tumour subtypes corresponded to
187 eight TCGA/ICGC cancer pathways (see Supplementary Methods, Extended Data
188 Fig. 5). While six showed slightly elevated mutational frequencies in African derived
189 tumours, genes impacting epigenetic mechanisms were significantly biased towards
190 Europeans (OR=0.179, *P*-value=2.9e-07, Extended Data Fig. 6b). Pathway
191 enrichment analysis supported five functional networks of the cancer pathways, with
192 two of them involved in signal transduction and DNA checkpoint processes to which
193 five of the eight pathways were interacted (Extended Data Fig. 6a; Supplementary
194 Table 8).

195 **Global molecular subtypes**

196 Through combining molecular profiling and patient demographics, ethnicity and
197 geography, we identify a new prostate cancer taxonomy we define as ‘Global
198 Mutational Subtypes (GMS)’ (Fig. 2b). While all European patients from Australia
199 (n=53) and Brazil (n=3) were limited to GMS-A and C, African derived tumours were
200 dispersed across all four subtypes. We found GMS-B and D to predominate in
201 Africans, with GMS-B including a single patient of admixed ancestry (92% African)
202 and GMS-D including a single admixed (63% African) and a single European
203 ancestral patient. The latter was one of only five Europeans in our study born and
204 raised in Africa. Compared to other patients of European ancestry, this patient showed
205 the highest mutational density across all types. Alternative consensus clustering of
206 individual mutational types mostly recapitulated the subtypes by integrative analysis
207 (Supplementary Table 6). Through further inclusion of Chinese Asian high-risk
208 prostate cancer data¹² (n=93, Extended Data Fig. 7a), we found GMS-A to be
209 ethnically and geographically ‘universal’, while GMS-D remained ‘African-specific’
210 with a new ‘African-Asian’ GMS-E emerging. GMS-B remained ‘African-specific’
211 and GMS-C ‘European-African’. While all patients were treatment naïve at the time
212 of sampling, our European cohort was recruited with extensive follow-up data
213 (median±s.d., 122.5±44.4 months). Interestingly, biochemical relapse (Fig. 3c) and
214 death-free survival probability (Fig. 3d) explains better clinical outcomes for patients
215 presenting with the ‘universal’ over the ‘European-African’ GMS (A *versus* C, log-
216 rank *P*-value=0.008 and 0.041, respectively).

217 Our novel GMS taxonomy could leverage pan-cancer studies in the following ways.
218 First, a sampling strategy of patients from the PCAWG project was rather
219 homogeneous in each cancer, therefore inhibiting the discovery of globally restricted
220 subtypes^{3,14} (Extended Data Fig. 7b). Second, ancestral²⁶ and geographic data of

221 patients should be included in molecular profiling of cancers. Lastly, the inclusion of
222 ethnic disparity in cancer studies would need to properly address admixture in a
223 sampling cohort, with too low ancestral cut-off appearing to create highly admixed,
224 but similar ancestry among individuals, therefore discouraging ethnically diverse
225 samples.

226 **Novel and known mutational signatures**

227 Approximating the contribution of mutational signatures to individual cancer
228 genomes facilitates an association of the signatures to exogenous or endogenous
229 mutagen exposures that contribute to the development of human cancer³. Here, we
230 generated a novel list of copy number (CN) and SV signatures and their contributions
231 to prostate cancer using nonnegative matrix factorisation²⁷ (Extended Data Fig. 8a, b).
232 Combined with a known catalogue of small mutational signatures, including single
233 base substitutions (SBS), doublet base substitutions (DBS) and small insertions and
234 deletions (ID), we observed not only a substantial variation in the number of
235 mutational features, but also over-representation in African derived tumours
236 (Extended Data Fig. 8c). Overall, 96 SBS, 78 DBS and 83 ID features examined had
237 significantly higher totals in Africans (SBS, 3,399 *versus* 2,840 in Europeans, *P*-
238 value=0.014; DBS, 42 *versus* 32, *P*-value=0.006; ID, 374 *versus* 360, *P*-value=0.016,
239 two-sample t-test). We generated six *de novo* signatures for each small signature type
240 (median cosine similarity 0.986, 0.856, and 0.976, respectively), corresponding to 12,
241 seven and eight global signatures, respectively (0.966, 0.850, and 0.946, respectively;
242 Extended Data Fig. 9), with 26 likely to be of biological origin (SBS47, possible
243 sequencing artefacts). DBS substitutions accounted for about 1% of the prevalence of
244 SBS. The CN features were also greater in Africans (CN, 3,971 *versus* 2,721, *P*-

245 value=1.92e-08; SV, 94 *versus* 88, *P*-value=0.100). The SV features defined in a
246 recent pan-cancer study²⁷ were each mutually exclusive and included simple SVs
247 (split according to size, replication timing and occurrence at fragile sites), templated
248 insertions (split by size), local n-jumps and local–distant clusters. The factorisation of
249 a sample-by-mutation spectrum matrix identified six CN signatures (CN1-6) and eight
250 SV signatures (SV1-8), as well as their contributions to each tumour.

251 We found the full spectrum of mutational signatures (SBS, DBS, ID, CN and SV) to
252 support our newly described GMS. Enrichment records of the top signatures in each
253 tumour were significantly associated type by type with the taxonomic subtypes,
254 except for DBS (*P*-values=5.1e-07–0.017, one-way ANOVA or Fisher’s exact test,
255 Extended Data Fig. 8d). Regardless of signature type, 13/40 mutational signatures
256 showed either inverse or proportionate correlations with our GMS (FDR=4.97e-13–
257 0.095, Spearman’s correlation, Fig. 4). Duplication signatures, including CN1
258 (tandem duplication), CN4 (whole genome duplication), SV2 (insertion) and SV5
259 (large duplication), were biased to the most mutationally noisy subtype (Extended
260 Data Fig. 8a, b), with CN4 and SV5 frequent in Africans (ρ =-0.24, FDR=0.005–
261 0.006). The mutational density of 30 out of 32 genes highly mutated in our GMS and
262 reported in prostate cancer was also significantly correlated with different somatic
263 signatures, with most observed in CN2, CN6 and SV6 signatures that were mainly
264 caused by deleted genomes. Small-size signatures were inversely significant among
265 20 mutated genes, indicating a higher number of mutations towards lesser mutated
266 tumours (FDR=1.05e-08–0.099).

267 **Life history of globally mutated subtypes**

268 Timeline estimates of individual somatic events reflect evolutionary periods that
269 differ from one patient to another; for example, a cluster of identical alterations
270 derived from clones in one patient presented as subclonal events in another patient
271 (Extended Data Fig. 10a, b). However, they provide in part the order of driver
272 mutations and CNAs present in each sample²⁵. The reconstruction of aggregating
273 single-sample ordering of all drivers and CNAs reveals different evolutionary patterns
274 unique to each GMS (Extended Data Fig. 10c, Fig. 5a, b). We draw approximate
275 cancer timelines for each GMS portraying the ordering of driver genes, recurrent
276 CNAs and signature activities chronologically interleaved with whole-genome
277 duplication (WGD) and the emergence of the most recent common ancestor (MRCA)
278 leading up to diagnosis. Basically, significantly co-occurring interactions of the
279 drivers and CNAs are shown (OR=2.6–97.8, P -values = 2.04e-30–0.01), supporting
280 their clonal and subclonal ordering states within the reconstructed timelines. SBS and
281 ID signatures that are abundant in each GMS display changes of their mutational
282 spectrum between the clonal and subclonal state, suggesting a difference in mutation
283 rates. The plot of clock-like CpG-to-TpG mutations and patient age adjustment shows
284 the median mutation rate as little as 0.968 per year for the ‘universal’ GMS, but the
285 highest rate at 1.315 per year observed in the ‘African-specific’ GMS-D. GMS-B and
286 C have rates of 1.144 and 1.092 per year, respectively. Assessing the relative timing
287 of somatic driver events, *TP53* mutations and accompanying 17p loss are of particular
288 interest, occurring early in GMS-C progression and at a later stage in GMS-A. League
289 model relative timing of driver events (see Supplementary Methods) is in line with a
290 fraction of probability distribution of the *TP53* alterations at the early stage, but most
291 are at an intermediate state of evolution (Extended Data Fig. 10d). This basic
292 knowledge of *in vivo* tumour development suggests that some tumours could have a

293 shorter latency period before reaching their malignant potential, so known genomic
294 heterogeneity of their primary clones is paramount to pave a way for early detection.

295 **Discussion**

296 To our knowledge, our study represents the first, if not, the largest whole-genome
297 prostate cancer, and likely any cancer, genome resource for Sub-Saharan Africa. Here
298 we describe a novel prostate cancer molecular taxonomy, identifying ethnically and
299 geographically distinctive Global Mutational Subtypes (GMS). Compared to previous
300 taxonomy using significantly mutated genes in prostate cancer^{7,19}, we found GMS to
301 compliment known subtypes such as *SPOP* and *FOXA1* mutations, in contrast to
302 underrepresented subtypes in this study, including gene fusions (Extended Data Fig.
303 4a). We also found GMS to correlate with mutational signatures reported in the
304 known catalogue of somatic mutations in cancer, where each tumour is represented by
305 different degrees of exogenous and endogenous mutagen exposures³. Our study has
306 leveraged the analysis of evolution across 38 cancer types by the PCAWG
307 Consortium²⁵, recognising that each GMS represents a unique evolutionary history
308 with drivers and mutational signatures varied between cancer stages and linking
309 somatic evolution to a patient's demographics. Therefore, some represent 'rare or
310 geographically restricted signatures' that are still a myth in pan-cancer studies^{3,14}.

311 We consider two extreme cases, 'universal' GMS-A *versus* 'African-specific' GMS-B
312 and D, that would have been influenced by two different mutational processes for
313 conceptual simplicity (Fig. 5c). One is predisposing genetic factors that are known for
314 prostate tumourigenesis across ethnolinguistic groups²⁸⁻³⁰. This factor contributes to
315 endogenous mutational processes, especially those with significant germline-somatic
316 interactions, such as the *TMPRSS2-ERG* fusion less frequently observed in men of

317 African and Asian ancestry^{12,31}, germline *BRCA2* mutations and the somatic *SPOP*
318 driver co-occurred with their respective counterparts^{32,33}. Another factor is modifiable
319 environmental attributes specific to certain circumstances or geographic regions that,
320 until now, have been elusive to prostate cancer. They act as mutagenic forces leading
321 to the positive selection of point mutations throughout life in healthy tissues^{34,35} and
322 cancers³⁶, forming fluid boundaries between normal ageing and cancer tissues.
323 According to Ottman³⁷, the above-mentioned model of gene-environment interaction
324 is observed when there is a different effect of a genotype on disease in individuals
325 with different environmental exposures or, alternatively, a different effect of an
326 environmental exposure on disease in individuals with different genotypes. Other
327 GMS subtypes would be a combination of the two processes, warranting a need for
328 larger populations of different ethnicities from different geographical localities to be
329 studied for a breakthrough in nature *versus* nurture. As such, the study directly
330 accounts for the large spatio-genomic heterogeneity of prostate cancer and its
331 associated evolutionary history in understanding the disease aetiology.

332 Our study suggests that larger genomic datasets of ethnically and geographically
333 diverse populations in a unified analysis will continue to identify rare and
334 geographically restricted subtypes in prostate cancer and potentially other cancers.
335 We are the first to demonstrate that ancestral and geographic attributes of patients
336 could facilitate those studies on cancer population genomics, an alternative to cancer
337 personalised genomics, for a better scientific understanding of nature *versus* nurture.

338 **Figure legends**

339 **Fig. 1 | Mutational density in prostate tumours of different ancestries. a**, Distribution of somatic
340 aberration (event number or number of base pairs) for seven mutational types across 183 tumour-blood
341 WGS pairs. **b**, Different types of mutational burden observed in this cohort. Samples are percentile
342 ranked and then ordered based on the sum of percentiles across the mutational types observed in each
343 ethnic group (left panel). Spearman's correlation is shown between mutation types, with dot size
344 representing the magnitude of correlation and background colour giving statistical significance of FDR
345 values (right panel).

346 **Fig. 2 | Prostate cancer taxonomy of ethnically diverse populations. a**, Integrative clustering
347 analysis reveals four distinct molecular subtypes of prostate cancer. The molecular subtypes are
348 illustrated by small somatic mutations (coding regions and noncoding elements), somatic copy number
349 alterations and somatic SVs. The percentage and association between the iCluster membership and
350 patient ancestry are illustrated in square brackets. A, African ancestry; Ad, Admixed; and E, European
351 ancestry. **b**, Total somatic mutations across four molecular subtypes in this study. Dashed lines indicate
352 the median values of mutational densities across the four subtypes. For each subtype, patients are
353 ordered based on their ethnicity.

354 **Fig. 3 | Aberration of driver genes in four diverse subtypes. a**, Analysis of the long tail of driver
355 genes using different mutation data combined. A total of 124 genes are associated with four prostate
356 cancer subtypes, and all have previously been reported as significantly recurrent mutations/SV
357 breakpoints in the PCAWG Consortium²², except for ones marked by asterisks, where they are
358 assigned to be significantly mutated using whole-genome data in this study. The Y-axis shows
359 corrected *P*-values in $-\log_{10} P$. CDS, coding driver data; NC, noncoding driver data; SV, significantly
360 recurrent breakpoint data; and CN, gene-level copy number data. **b**, Unsupervised hierarchical
361 clustering of known and putative driver genes identified within four prostate cancer subtypes (A-D, a
362 bottom-up direction). Rows are patients, and columns represent 124 driver genes (alphabetical order)
363 identified using different mutational types. **c**, Kaplan-Meier plot of biochemical relapse (BCR)-free
364 survival proportion of European patients in subtype A (n=161) *versus* C (n=19). **d**, Kaplan-Meier plot
365 of cancer survival probability of European patients in subtype A (n=82) *versus* C (n=17).

366 **Fig. 4 | Estimates of genomic aberrations contributed by each mutational signature.** The size of
367 each dot represents FDR values of Spearman correlation P -values using BH correction. The colours of
368 each dot represent correlation coefficient (ρ). GMS is assigned as 1-4 for Subtypes A-D,
369 respectively; African, Admixed and European are recorded as 1-3, respectively. The correlation of 32
370 significantly mutated genes in prostate cancer is shown in the X-axis.

371 **Fig. 5 | Evolutionary history of globally mutated subtypes. a,** The cancer timeline of the universal
372 subtype begins from the fertilised egg to the age of the patients at a cohort. **b,** that of GMS-C.
373 Estimates for major events, such as WGD (whole-genome duplication) and the emergence of the
374 MRCA (the most recent common ancestor), are used to define early, variable, late and subclonal stages
375 of tumour evolution approximately in chronological time. When early and late clonal stages are
376 uncertain, the variable stage is assigned. The variable/constant time period includes events that are
377 ranked before the WGD event and also begins shortly after another break in the timeline. The late
378 period does have a definite start, as this includes events that are ranked after WGD, when it occurs.
379 Driver genes and CNAs are shown in each stage if present in previous studies^{8,22} and defined by
380 MutationTime.R program. Mutational signatures (Sigs) that, on average, change over the course of
381 tumour evolution, or are substantially active but not changing, are shown in the epoch in which their
382 activity is rather greatest. Dagger symbols denote alterations that are found to have different timing.
383 Significant pairwise interaction events between the mutations and copy number alterations were
384 computed using Odds Ratio (OR). Either co-occurrence or mutually exclusive event is considered if
385 $OR > 2$ or < 0.5 , respectively. Median mutation rates of CpG-to-TpG burden per Gb are calculated using
386 age-adjusted branch length of cancer clones and maximally branching subclones. **c,** Schematic
387 representation of a world map with the distribution of GMS (A–D) among ethnically/globally diverse
388 populations. The gene-environment interaction model of globally mutated subtypes is shown in the
389 right panel. The contingency table of number of patients with different ancestries (germline variants)
390 stratified by subtypes and associated with certain geography or environmental exposure (two-sided P -
391 value= 0.0005, Fisher's exact test with 2,000 bootstraps).

392 **Methods**

393 **Patient cohorts and whole-genome sequencing**

394 Our study included ~180 treatment naïve prostate cancer patients recruited under
395 informed consent and appropriate ethics approval (Supplementary Methods, Section
396 2) from Australia (n=53), Brazil (n=7) and South Africa (n=123). DNA extracted
397 from fresh tissue and matched blood underwent 2x150 bp sequencing on the Illumina
398 NovaSeq instrument (Kinghorn Centre for Clinical Genomics, Garvan Institute of
399 Medical Research).

400 **WGS processing and variant calling**

401 Each lane of raw sequencing reads was aligned against human reference hg38 +
402 alternate contigs using bwa v0.7.15³⁸. Lane-level BAMs from the same library were
403 merged, and duplicate reads were marked. The Genome Analysis Toolkit (GATK
404 v4.1.2.0) was used for base quality recalibration³⁹. Contaminated and duplicate
405 samples (n=8) were removed. We implemented three main pipelines for the discovery
406 of germline and somatic variants, with the latter including small (SNV and indel) to
407 large genomic variation (CN and SV). Complete pipelines and tools used are available
408 from the Sydney Informatics Hub (SIH), Core Research Facilities, University of
409 Sydney (see Code availability). Scalable bioinformatic workflows are described in
410 Supplementary Methods, Section 4.

411 Genetic ancestry was estimated using fastSTRUCTURE v1.0⁴⁰, Bayesian inference
412 for the best approximation of marginal likelihood of a very large variant dataset.
413 Reference panels for African and European ancestry compared in this study were
414 retrieved from previous whole-genome databases^{20,21}.

415 **Analysis of chromothripsis and chromoplexy**

416 Clustered genomic rearrangements of prostate tumours were identified using
417 ShatterSeek v0.4⁴¹ and ChainFinder v1.0.1⁴². Our somatic SV and somatic CNA
418 callsets were prepared and co-analysed using custom scripts (see Code availability,
419 Supplementary Methods, Section 6).

420 **Analysis of mutational recurrence**

421 We used three approaches to detect recurrently mutated genes or regions based on
422 three mutational types, including small mutations, SVs and CNAs (see Supplementary
423 Methods, Section 7). In brief, small mutations were tested within a given genomic
424 element as being significantly more mutated than adjacent background sequences.
425 The genomic elements retrieved from syn5259886, the PCAWG Consortium²² were a
426 group of coding sequences and 10 groups of noncoding regions. SV breakpoints were
427 tested in a given gene for their statistical enrichment using Gamma-Poisson regression
428 and corrected by genomic covariates¹³. Focal and arm-level recurrent CNAs were
429 examined using GISTIC v2.0.23⁴³. Known driver mutations in coding and noncoding
430 regions published in PCAWG^{22,44,45} were additionally recorded in our 183 tumours,
431 and those specific to prostate cancer genes were also included^{7,8,13,18,19}.

432 **Integrative analysis of prostate cancer subtypes**

433 Integrative clustering of three genomic data types for 183 patients was performed
434 using iClusterplus^{12,46} in R, with the following inputs: *i*) driver genes and elements; *ii*)
435 somatic CN segments; and *iii*) significantly recurrent SV breakpoints. We ran
436 iClusterPlus.tune with clusters ranging from 1-9. We also performed unsupervised
437 consensus clustering on each of the three data types individually. Association analysis

438 of genomic alteration with different iCluster subtypes was performed in detail in
439 Supplementary Methods, Section 8. Differences in drivers, recurrent breakpoints and
440 somatic CNAs across different iCluster subtypes were reported.

441 **Comparison of iCluster with Asian and pan-cancer data**

442 To compare molecular subtypes between extant human populations, the Chinese
443 Prostate Cancer Genome and Epigenome Atlas (CPGEA, PRJCA001124)¹² was
444 merged and processed with our integrative clustering analysis across the three data
445 types described above, with some modifications. Moreover, we leveraged the
446 PCAWG Consortium¹⁴ to define molecular subtypes across different ethnic groups in
447 other cancer types using published data of somatic mutations, SV and GISTIC results
448 by gene. Four cancer types that consisted of breast, liver, ovarian, and pancreatic
449 cancers were considered due to existing primary ancestries of African, Asian and
450 European with at least 70% contribution. Full details are given in Supplementary
451 Methods, section 8.4.

452 Prostate cancer subjects of PCAWG¹⁴ were retrieved to compare with Australian data
453 with clinical follow-up. Only those of European ancestry greater than 90% (n=139)
454 were analysed for the three genomic data types of iCluster subtyping, as well as
455 individual consensus clustering. Clustering results identical to the larger cohort size
456 mentioned above were chosen for association analyses. Differences in the
457 biochemical relapse and lethal prostate cancer of the subjects across the subtypes
458 were assessed using the Kaplan–Meier plot followed by a log-rank test for
459 significance.

460 **Analysis of mutational signatures**

461 Mutational signatures (SBS, DBS and ID), as defined by the PCAWG Mutational
462 Signatures Working Group³, were fit to individual tumours with observed signature
463 activity using SigProfiler⁴⁷. Nonnegative matrix factorisation (NMF) was
464 implemented to detect *de novo* and global signature profiles among 183 patients and
465 their contributions. Novel mutational genome rearrangement signatures (CN and SV)
466 were also performed using the NMF, with 45 CN and 44 SV features examined across
467 183 tumours. We followed the PCAWG working classification and annotation scheme
468 for genomic rearrangement²⁷. Two SV callers were used to obtain exact breakpoint
469 coordinates. Replication timing scores influencing on SV detection were set at >75,
470 20-75, and <20 for early, mid, and late timing, respectively⁴⁸. Full details of analysis
471 steps, parameters and relevant statistical tests are given in Supplementary Methods,
472 Section 9.

473 **Reconstruction of cancer timelines**

474 Timing of copy number gains and driver mutations (SNVs and indels) into four
475 epochs of cancer evolution (early clonal, unspecified clonal, late clonal, and
476 subclonal) was conducted using MutationTimeR²⁵. CN gains including 2+0, 2+1, and
477 2+2 (1+1 for a diploid genome) were considered for a clearer boundary between
478 epochs instead of solely information of variant allele frequency. Confidence intervals
479 ($t_{lo} - t_{up}$) for timing estimates were calculated with 200 bootstraps. Mutation rates for
480 each subtype were calculated following Gerstung, et al²⁵ that CpG-to-TpG mutations
481 were counted for the analysis because they were attributed to spontaneous
482 deamination of 5-methyl-cytosine to thymine at CpG dinucleotides, therefore acting
483 as a molecular clock.

484 League model relative ordering was performed to aggregate across all study samples
485 to calculate the overall ranking of driver mutations and recurrent CNAs. The
486 information for the ranking was derived from the timing of each driver mutation and
487 that of clonal and subclonal CN segments, as described above. Full description is
488 provided in Supplementary Methods, Section 10.

489 **Data availability**

490 Alignments, somatic and germline variant calls, annotations and derived datasets are
491 available for general research use for browsing and download through the European
492 Genome-Phenome Archive (accession number [EGA0000000000](#)). Other supporting
493 data are available upon request from the corresponding author.

494 **Code availability**

495 The core computational pipelines used in this study for read alignment, quality control
496 and variant calling are available to the public at [https://github.com/Sydney-](https://github.com/Sydney-Informatics-Hub/Bioinformatics)
497 [Informatics-Hub/Bioinformatics](https://github.com/Sydney-Informatics-Hub/Bioinformatics). Analysis code for chromothripsis and chromoplexy
498 is available through another GitHub page, https://github.com/tgong1/Code_HRPCa.

499 **Acknowledgements**

500 The work presented was supported by the National Health and Medical Research
501 Council (NHMRC) of Australia through a Project Grant (APP1165762, V.M.H.),
502 NHMRC Ideas Grant (APP2001098, V.M.H. and M.S.R.B.), University of Sydney
503 Bridging Grant (G199756, V.M.H.), and partly through the U.S. Department of
504 Defense (DoD) Prostate Cancer Research Program (PCRP) Idea Development Award
505 (PC200390, including W.J., S.M.P., D.C.W., S.M., M.S.R.B. and V.M.H.). The

506 authors acknowledge the use of the National Computational Infrastructure (NCI)
507 which is supported by the Australian Government, and accessed through the National
508 Computational Merit Allocation Scheme (V.M.H., E.K.F.C and W.J.), the Intersect
509 Computational Merit Allocation Scheme (V.M.H.), Intersect Australia Limited, and
510 the Sydney Informatics Hub, Core Research Facility, while we acknowledge the
511 Garvan Institute of Medical Research's Kinghorn Centre for Clinical Genomics
512 (KCCG) core facility for data generation. Recruitment, sampling and processing for
513 the Southern African Prostate Cancer Study (SAPCS), as required for the purpose of
514 this study, was supported by the Cancer Association of South Africa (CANSA,
515 M.S.R.B. and V.M.H.). V.M.H. was supported by Petre Foundation via the University
516 of Sydney Foundation, A-M.H. and W.J. by a Cancer Institute of New South Wales
517 (CINSW) Program Grant (TPG172146 to L.G.H., J.G.K., P.D.S. and V.M.H.), with
518 additional support to W.J. provided by the Prostate Cancer Research Alliance
519 Australian Government and Movember Foundation Collaboration PRECEPT
520 (Prostate cancer prognosis and treatment study, led by A/Prof. N. Corcoran,
521 University of Melbourne, Australia). T.G. is now located at the Human Phenome
522 Institute, Fudan University, Shanghai, China and E.K.F.C. at NSW Health Pathology,
523 Sydney, Australia. We are forever grateful to the patients and their families who have
524 contributed to this study; without their contribution, this research would not be
525 possible. We acknowledge the contributions of the many clinical staff across the
526 SAPCS (South Africa), the St Vincent's Hospital Sydney (Australia) and
527 6LEndocrine and Tumor Molecular Biology Laboratory (Brazil), who over many
528 years have recruited patients and provided samples to these critical bioresources, with
529 special recognition of Professor Philip Venter (retired), Dr's Richard L. Monare

530 (retired) and Dr Smit van Zyl, previously from the University of Limpopo, South
531 Africa, for their critical contributions as inaugural members of the SAPCS.

532 **Authors' contributions**

533 V.M.H. designed the experiments and supervised the project; W.J. led the
534 bioinformatic and statistical analyses, while both W.J. and V.M.H. performed data
535 interpretation. S.M.P., R.J.L., A-M.H., and D.G.P. prepared the samples and managed
536 phenotypic data. M.L. and J.G.K. performed pathological grading, while R.C.,
537 L.G.H., I.S.B., S.B.A.M., P.D.S. and M.S.R.B. managed patient recruitments and
538 consents, as well as clinical interpretation. V.M.H., S.B.A.M. and M.S.R.B. codirect
539 the Southern African Prostate Cancer Study (SAPCS). W.J., J.J., T.G., C.W., T.C. and
540 R.S. developed the pipelines and performed the efficient and scalable high-
541 performance computational variant calling, with critical advice provided by E.K.F.C
542 and V.M.H. W.J., J.J. and T.G. performed complex variant annotation, while R.J.L.
543 generated the optical genome mapping (OGM) data. W.J. performed mutational
544 signature and tumour evolution analysis, with critical advice provided by D.C.W.
545 W.J. and V.M.H. wrote the manuscript. W.J. generated the figures, while all authors
546 contributed to the final editing and approval.

547 **Competing interest declaration**

548 The authors declare no competing interests.

549 **Supplementary Tables**

550 **Supplementary Table 1 | Clinical cohort characteristics and sequencing quality**

551 **Supplementary Table 2 | Driver information by patient**

552 **Supplementary Table 3 | GISTIC2 results of all genomic lesions under 99% confidence level**

553 **Supplementary Table 4 | List of significantly recurrent SV breakpoints at FDR lower than 0.10**

554 **Supplementary Table 5 | TCGA prostate cancer taxonomy identified in this study**

555 Patient by driver mutation and patient by driver structural variation summary matrices are provided.

556 **Supplementary Table 6 | Integrative iCluster analysis of 183 prostate tumours**

557 **Supplementary Table 7 | List of 124 preferentially mutated genes within four tumour subtypes**

558 **Supplementary Table 8 | Pathway enrichment analysis of 124 preferentially mutated genes**

559 **Supplementary Table 9 | Total mutational signature profiles across 183 tumours**

560 The table shows data matrices of SBS feature by patient, DBS feature by patient, ID feature by patient,

561 CN feature by patient, and SV feature by patient.

562 **Supplementary Table 10 | Cross-individual contamination level**

563 **Supplementary Table 11 | Cancer evolution analysis of prostate cancer**

564 Clonal architecture by PhyloWGS and timing of gains and drivers by MutationTimeR is provided per

565 tumour

566

567

568 **Extended data legends**

569 **Extended Data Fig. 1 | Clinical cohorts and statistical metrics.** **a**, Clinical and pathological patient
570 characterisation. **b**, STRUCTURE analysis of bi-allelic germline variants with the logistic prior model.
571 Model components used to explain structure in the plot are $K=5$. All spectrum of African contributions
572 are summed and assigned as African ancestry. **c**, Saturation curve for all driver types across 183
573 patients. Recurrent copy number gains and losses were measured using GISTIC v2 (Supplementary
574 Methods). CDS, coding sequence; SV, structural variation. **d**, Spearman's correlation between different
575 variables measured in this cohort. Dot sizes represent the magnitude of correlation, with significant P -
576 values <0.01 .

577

578 **Extended Data Fig. 2 | Somatic driver mutations in 183 prostate cancer patients.** The covariates
579 on the left show mutational types and statistical significance (FDR) from ActiveDriverWGS and
580 GISTIC2. **a**, The top 300 driver genes in PCAWG discovered in primary prostate tumours among 183
581 specimens. The top barplot shows the distribution of the number of prostate cancer drivers and/or that
582 of PCAWG. The heatmap shows drivers found in this study (rows) for each patient (columns).
583 Heatmaps are coloured by mutational type. Bottom covariates show the clinical features of patients.
584 The percentage of transition/transversion mutations across 183 patients shows 1,364,210 small somatic
585 mutations across chromosomes 1-Y. **b**, The bottom heatmap shows the top 75 of previously reported
586 coding driver genes in prostate cancer observed in this study^{7,8,18,19}. The right barplot shows the number
587 of patients for each driver.

588

589 **Extended Data Fig. 3 | Discovery of prostate cancer drivers.** **a**, The number and types of PCAWG
590 driver genes and elements studied in our cohort. **b**, Recurrent copy number alterations among 183
591 prostate tumours identified with a 99% confidence level using GISTIC v2 (Supplementary Methods).
592 The figure shows GISTIC peaks of significant regions of recurrent amplification (red) or deletion
593 (blue) supported by $FDR <0.01$. **c**, Genome-wide scan for significantly recurrent breakpoints in our
594 study. The quantile-quantile plot shows P -values for mutational densities across 183 prostate cancer
595 patients. Generalised linear modelling (GLM) of somatic mutation densities along the genome with
596 significant background mutational processes adjusted in the model is also shown. **d**, Bionano
597 Genomics optical genome mapping at the HLA complex. Examples of HLA translocations from a

598 European patient (ID 12543) and an African patient (ID UP2360) studied in this cohort are
599 characterised by pairs of optical maps, each carrying a fusion junction with flanking fragments aligning
600 to one side of the two reference breakpoints. Using the recurrent HLA breakpoints identified in this
601 study, the genome map of the African specimen is found to have a low-end fusion junction matched
602 with chromosome 6 through a manual inspection of unfiltered consensus maps using Bionano Access
603 v15.2. Note that the HLA alternate contig fused in the European tumour is different from one suggested
604 by short-read sequencing (chr6_GL000252v2_alt). The reference genome map is an *in silico* digest of
605 the human reference hg38 with the DLE-1 enzyme. Genome map sizes are indicated on the horizontal
606 axis in megabase (Mb) units. Matching fluorescent labels between sample and reference genome map
607 are connected by gray lines.

608 **Extended Data Fig. 4 | TCGA molecular taxonomy.** **a**, Seven important oncogenic drivers identified
609 by TCGA within our African and European patients. **b**, Coding mutations observed within *SPOP* and
610 *FOXAI* genes. Rarely, a mutation at the BTB domain of *SPOP* gene is shown (R221C in an African
611 patient, KAL0072). FH, forkhead. **c**, *ETVI* fusions within positive patients caused by copy number
612 (CN) losses and/or structural variants (DEL, deletion; ICX, interchromosomal translocation; and INV,
613 unbalanced or balanced inversion). CN changes in chromosome 7 show the *ETVI* loss with log₂ CN
614 ratio less than -0.2. **d**, *ERG* fusions caused by CN losses and/or structural variants.

615

616 **Extended Data Fig. 5 | Prostate cancer genes and pathways.** The search is carried out using the
617 TCGA and ICGC cancer databases. The top affected genes for each pathway are present with lollipop
618 plots to show their hotspots of simple coding mutations if they existed.

619 **Extended Data Fig. 6 | Major biological pathways and networks of prostate cancer.** **a**, Networks
620 of functional interactions between driver genes are shown for each cancer pathway. Nodes represent
621 Gene Ontology biological processes and Reactome pathways and edges show functional interactions.
622 **b**, Pathway alteration frequencies between African and European. A sample was considered altered in
623 a given pathway if at least a single gene in the pathway had a genomic alteration. *P*-values indicate the
624 level of significance (two-sided Fisher's exact test).

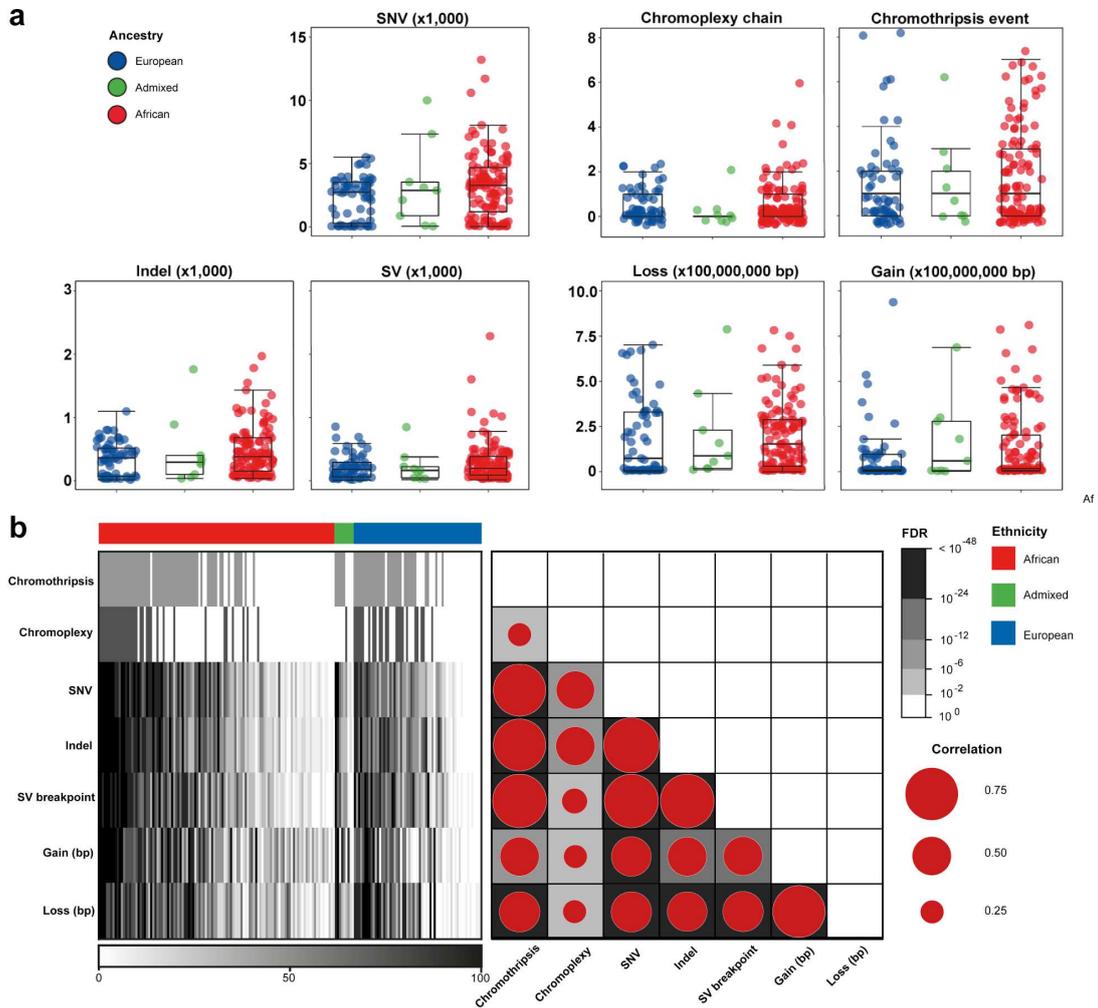
625 **Extended Data Fig. 7 | Molecular subtypes in prostate cancer and pan-cancers.** **a**, Unsupervised
626 hierarchical clustering of primary prostate tumours across three major ethnic groups was performed
627 using total somatic mutations present within WGS normalised data. Admixed individuals were also
628 tested in prostate cancer subtypes to which they belonged. **b**, Molecular subtyping of total somatic
629 mutations within pan-cancer studies, namely pancreatic, ovarian, breast and liver cancers. Raw data of
630 small somatic mutations, structural variants and copy number alterations acquired per cancer were
631 retrieved from the PCAWG¹⁴. For each subtype, patients are ordered based on their ethnicity. Ethnic
632 groups are assigned using a cut-off of ancestral contribution greater than 70%; otherwise, considered as
633 Admixed.

634 **Extended Data Fig. 8 | Known and novel mutational signatures in prostate cancer.** **a**, Copy
635 number signatures in prostate cancer across 45 CN features ranked by mutational processes observed.
636 The six most distinctive signatures and their important components extracted by the NMF algorithm
637 were run on the sample size of 183 genomes. Bar charts represent the estimated proportion of each
638 event feature assigned to each signature (rows sum to one). **b**, Structural variation signatures in prostate
639 cancer ranked by mutational processes observed from small deletion to reciprocal rearrangement. The
640 eight most distinctive signatures and their important components extracted from 44 features using the
641 NMF algorithm were run on the sample size of 183 genomes. Bar charts represent the estimated
642 proportion of each event feature assigned to each signature (rows sum to one). **c**, Frequency of SBS,
643 DBS, ID, CN and SV features across 183 tumours. Colours at the bottom panel show the following
644 ethnic groups: *i*) African, red; *ii*) Admixed, green; and *iii*) European, blue. **d**, Stacked barplots of
645 multiple signature exposures for each mutational type enriched per patient and ranked by ethnic group.
646 Copy number and structural variation signatures (CN1-6 and SV1-8, respectively) are the first
647 identified in this study for prostate cancer, and their enrichment in a patient appears to be significantly
648 associated (P -values <0.05) with our GMS, considering either *de novo* or global mutational signatures
649 discovered in the Catalogue of Somatic Mutations in Cancer (COSMIC).

650 **Extended Data Fig. 9 | Total profiles of SBS, DBS, ID, CN and SV signatures.** The classification of
651 each signature type (SBS, 96 classes; DBS, 78 classes; ID, 83 classes; CN, 45 classes; and SV, 44

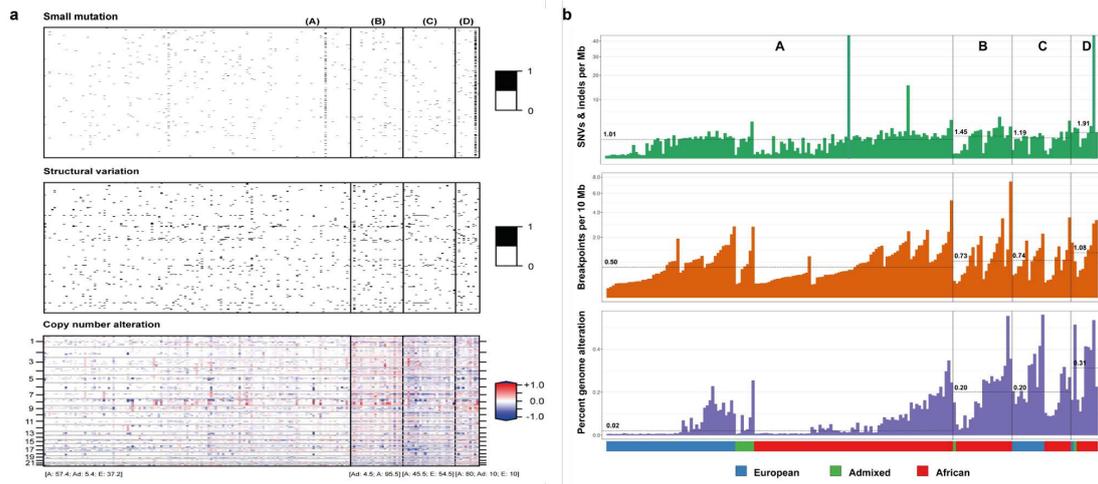
652 classes) is described in Supplementary Methods. The plotted data are available in digital form
653 (Supplementary Table 9).

654 **Extended Data Fig. 10 | Stages of prostate tumour development. a,** Clonal architecture and its
655 frequency in prostate cancer between Africans and Europeans. Tumours are divided into three groups:
656 monoclonal, linear and branching polyclonal. The number of small somatic mutations (SSM) and CNA
657 as percentage of genome alteration (PGA) is provided as median and range in bracket. Cancer cell
658 fraction (CCF) in each clone and/or subclone is shown in a circular node. Tumours that show
659 characteristics consistent with being polytumours or with multiple independent primary tumors are
660 excluded to remain conservative. **b,** Unbiased hierarchical clustering of CNA between clonal (trunk)
661 and subclonal (branch) mutations. Trunk mutations encompass those that occur between the root node
662 (normal) and its only child node, while all others are classified to have occurred in branch. Red
663 indicates gain; blue indicates loss; and rows indicate patients. Unidentified regions in trunk and branch
664 are assumed to have neutral copy number. ConsensusClusterPlus showed seven CNA clusters among
665 our patients to be optimal. The figure shows that a trunk alteration from one patient is mutationally
666 similar to a branch alteration from another, rather than to other trunk ones from different patients in a
667 cohort. **c,** Cancer timelines of GMS-B and D identified in this study. Detailed explanation is provided
668 in Fig. 5. **d,** Relative ordering model (PhylogicNDT LeagueModel) results for a cohort of samples
669 (n=66). The samples can be analysed if they have somatic events of interest prevalent greater than 5%
670 of the sample size and have informative clonal status available for each event (16 events). Probability
671 distributions show the uncertainty of timing for specific events in the cohort.



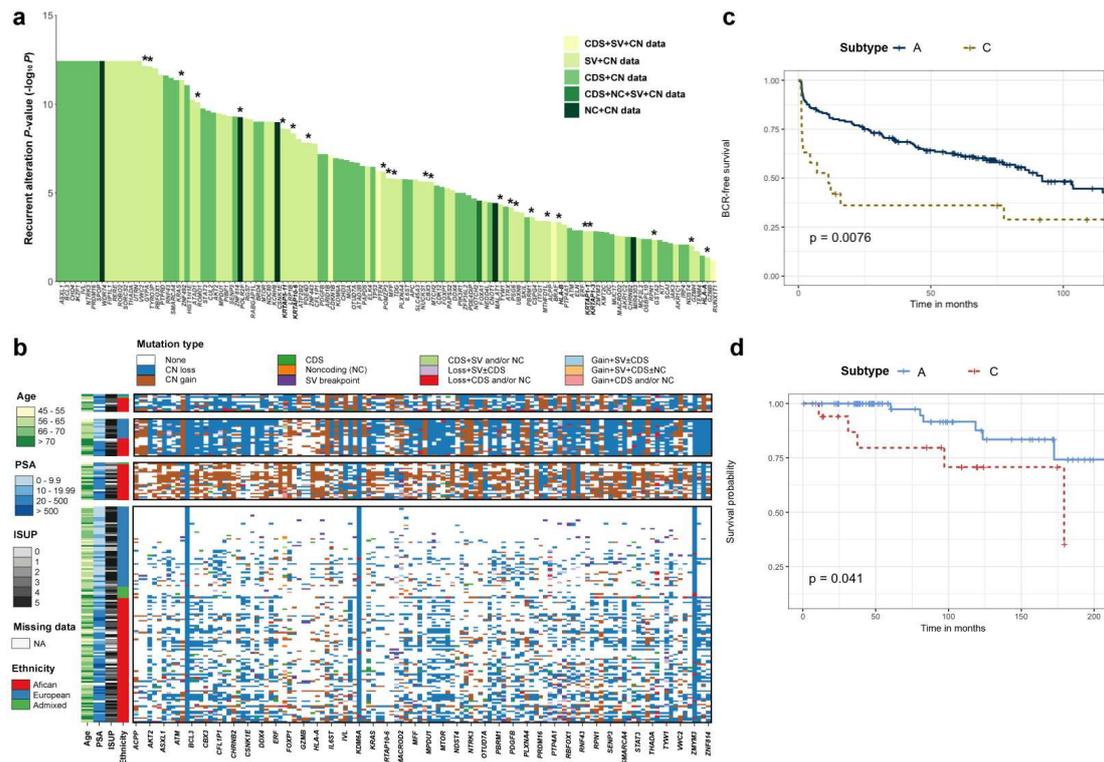
673

674 **Fig. 1 | Mutational density in prostate tumours of different ancestries. a**, Distribution of somatic
 675 aberration (event number or number of base pairs) for seven mutational types across 183 tumour-blood
 676 WGS pairs. **b**, Different types of mutational burden observed in this cohort. Samples are percentile
 677 ranked and then ordered based on the sum of percentiles across the mutational types observed in each
 678 ethnic group (left panel). Spearman’s correlation is shown between mutation types, with dot size
 679 representing the magnitude of correlation and background colour giving statistical significance of FDR
 680 values (right panel).



681

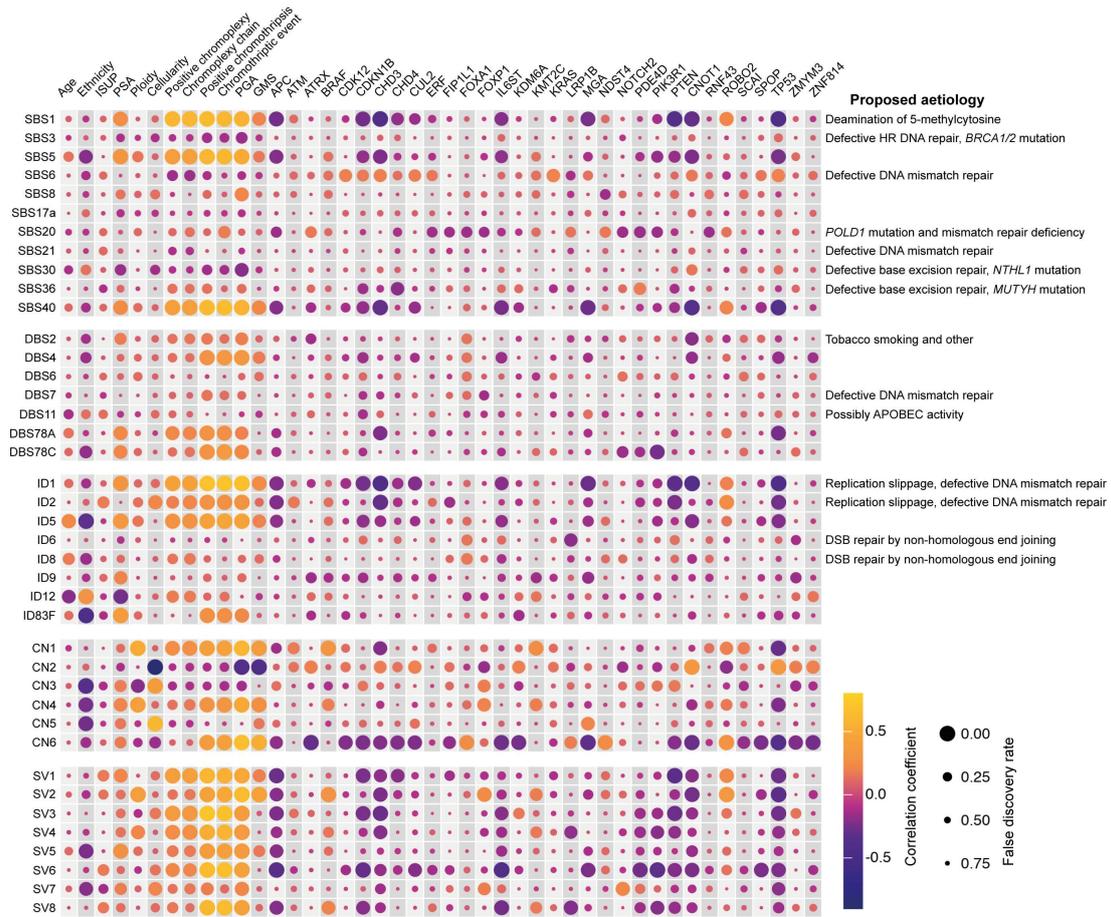
682 **Fig. 2 | Prostate cancer taxonomy of ethnically diverse populations. a**, Integrative clustering
 683 analysis reveals four distinct molecular subtypes of prostate cancer. The molecular subtypes are
 684 illustrated by small somatic mutations (coding regions and noncoding elements), somatic copy number
 685 alterations and somatic SVs. The percentage and association between the iCluster membership and
 686 patient ancestry are illustrated in square brackets. A, African ancestry; Ad, Admixed; and E, European
 687 ancestry. **b**, Total somatic mutations across four molecular subtypes in this study. Dashed lines indicate
 688 the median values of mutational densities across the four subtypes. For each subtype, patients are
 689 ordered based on their ethnicity.



690

691 **Fig. 3 | Aberration of driver genes in four diverse subtypes.** **a**, Analysis of the long tail of driver
 692 genes using different mutation data combined. A total of 124 genes are associated with four prostate
 693 cancer subtypes, and all have previously been reported as significantly recurrent mutations/SV
 694 breakpoints in the PCAWG Consortium²², except for ones marked by asterisks, where they are
 695 assigned to be significantly mutated using whole-genome data in this study. The Y-axis shows
 696 corrected P -values in $-\log_{10}P$. CDS, coding driver data; NC, noncoding driver data; SV, significantly
 697 recurrent breakpoint data; and CN, gene-level copy number data. **b**, Unsupervised hierarchical
 698 clustering of known and putative driver genes identified within four prostate cancer subtypes (A-D, a
 699 bottom-up direction). Rows are patients, and columns represent 124 driver genes (alphabetical order)
 700 identified using different mutational types. **c**, Kaplan-Meier plot of biochemical relapse (BCR)-free
 701 survival proportion of European patients in subtype A ($n=161$) *versus* C ($n=19$). **d**, Kaplan-Meier plot
 702 of cancer survival probability of European patients in subtype A ($n=82$) *versus* C ($n=17$).

703

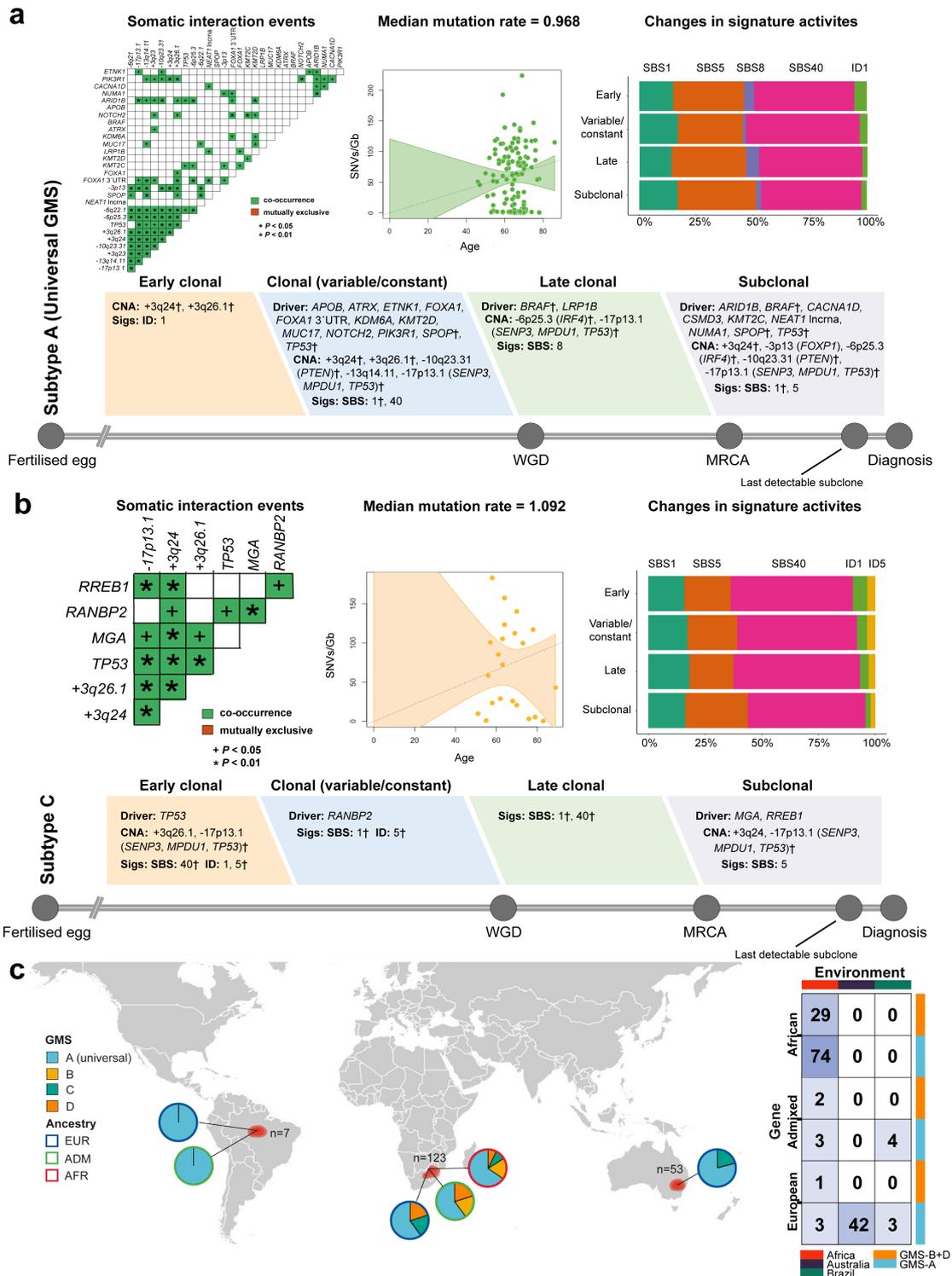


704

705 **Fig. 4 | Estimates of genomic aberrations contributed by each mutational signature.** The size of
 706 each dot represents FDR values of Spearman correlation P -values using BH correction. The colours of
 707 each dot represent correlation coefficient (ρ). GMS is assigned as 1-4 for Subtypes A-D,
 708 respectively; African, Admixed and European are recorded as 1-3, respectively. The correlation of 32
 709 significantly mutated genes in prostate cancer is shown in the X-axis.

710

711



712

713 **Fig. 5 | Evolutionary history of globally mutated subtypes. a**, The cancer timeline of the universal

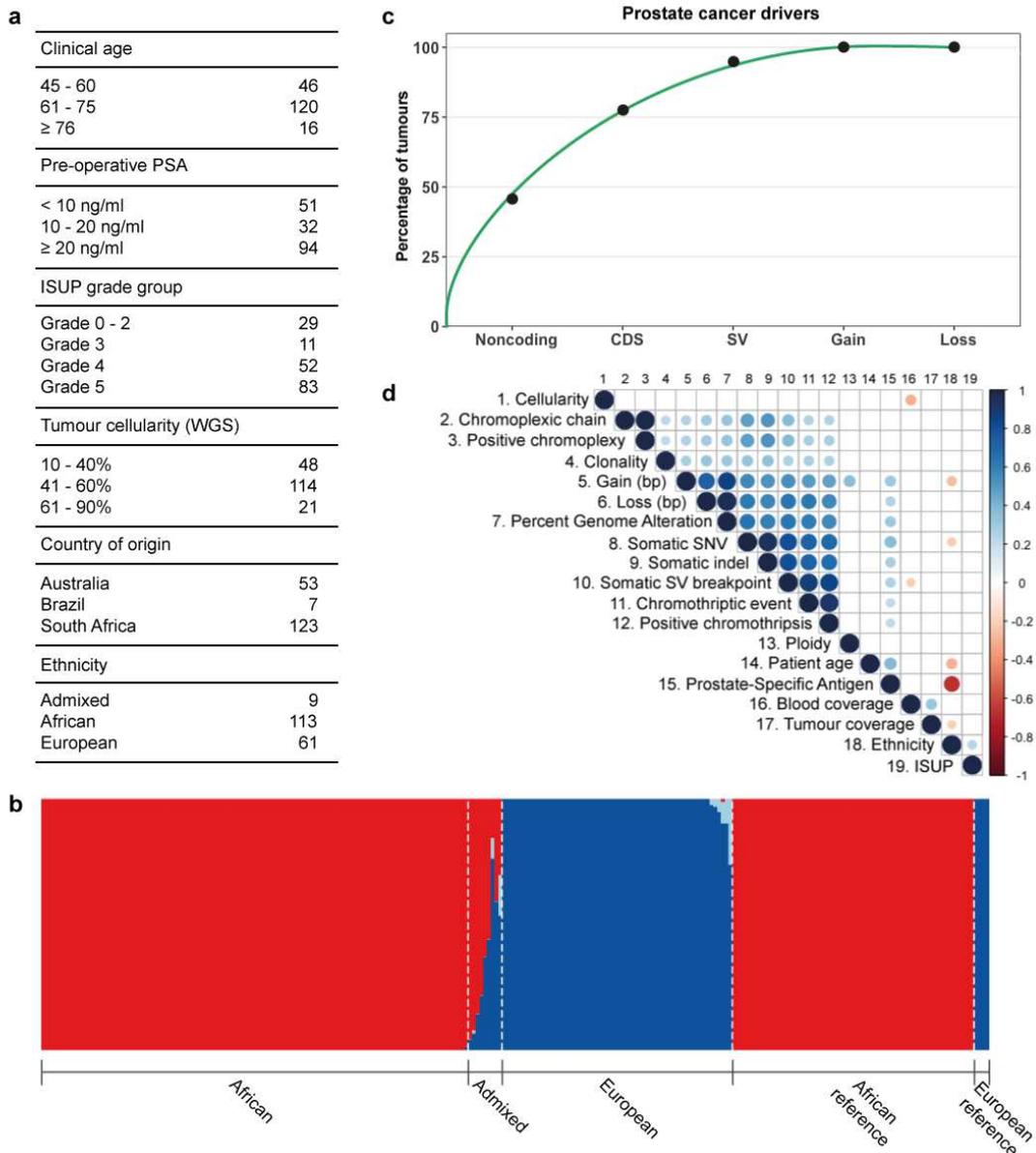
714 subtype begins from the fertilised egg to the age of the patients at a cohort. **b**, that of GMS-C.

715 Estimates for major events, such as WGD (whole-genome duplication) and the emergence of the

716 MRCA (the most recent common ancestor), are used to define early, variable, late and subclonal stages

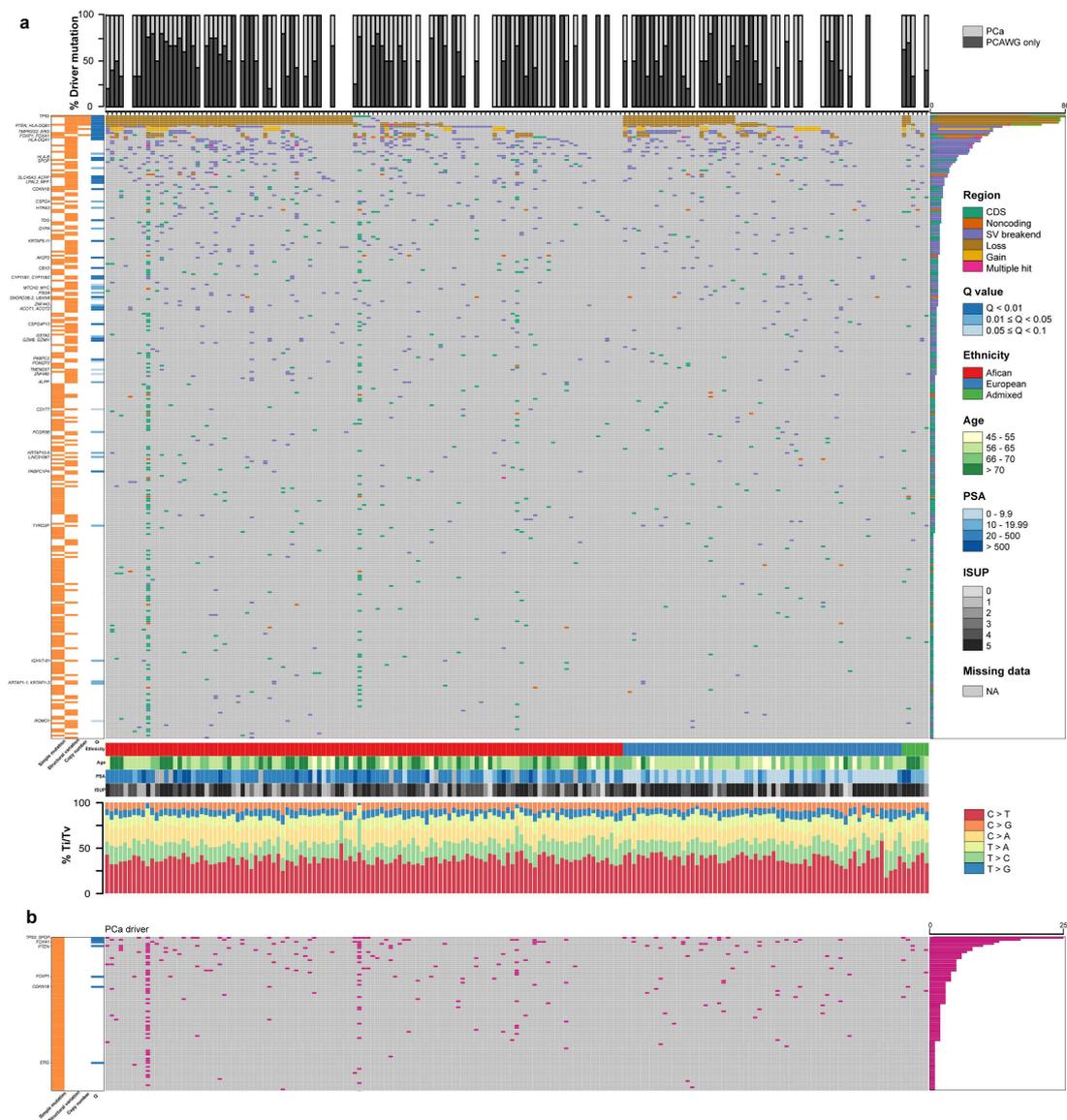
717 of tumour evolution approximately in chronological time. When early and late clonal stages are
718 uncertain, the variable stage is assigned. The variable/constant time period includes events that are
719 ranked before the WGD event and also begins shortly after another break in the timeline. The late
720 period does have a definite start, as this includes events that are ranked after WGD, when it occurs.
721 Driver genes and CNAs are shown in each stage if present in previous studies^{8,22} and defined by
722 MutationTime.R program. Mutational signatures (Sigs) that, on average, change over the course of
723 tumour evolution, or are substantially active but not changing, are shown in the epoch in which their
724 activity is rather greatest. Dagger symbols denote alterations that are found to have different timing.
725 Significant pairwise interaction events between the mutations and copy number alterations were
726 computed using Odds Ratio (OR). Either co-occurrence or mutually exclusive event is considered if
727 $OR > 2$ or < 0.5 , respectively. Median mutation rates of CpG-to-TpG burden per Gb are calculated using
728 age-adjusted branch length of cancer clones and maximally branching subclones. **c**, Schematic
729 representation of a world map with the distribution of GMS (A–D) among ethnically/globally diverse
730 populations. The gene-environment interaction model of globally mutated subtypes is shown in the
731 right panel. The contingency table of number of patients with different ancestries (germline variants)
732 stratified by subtypes and associated with certain geography or environmental exposure (two-sided *P*-
733 value= 0.0005, Fisher's exact test with 2,000 bootstraps).

734 **Extended data**



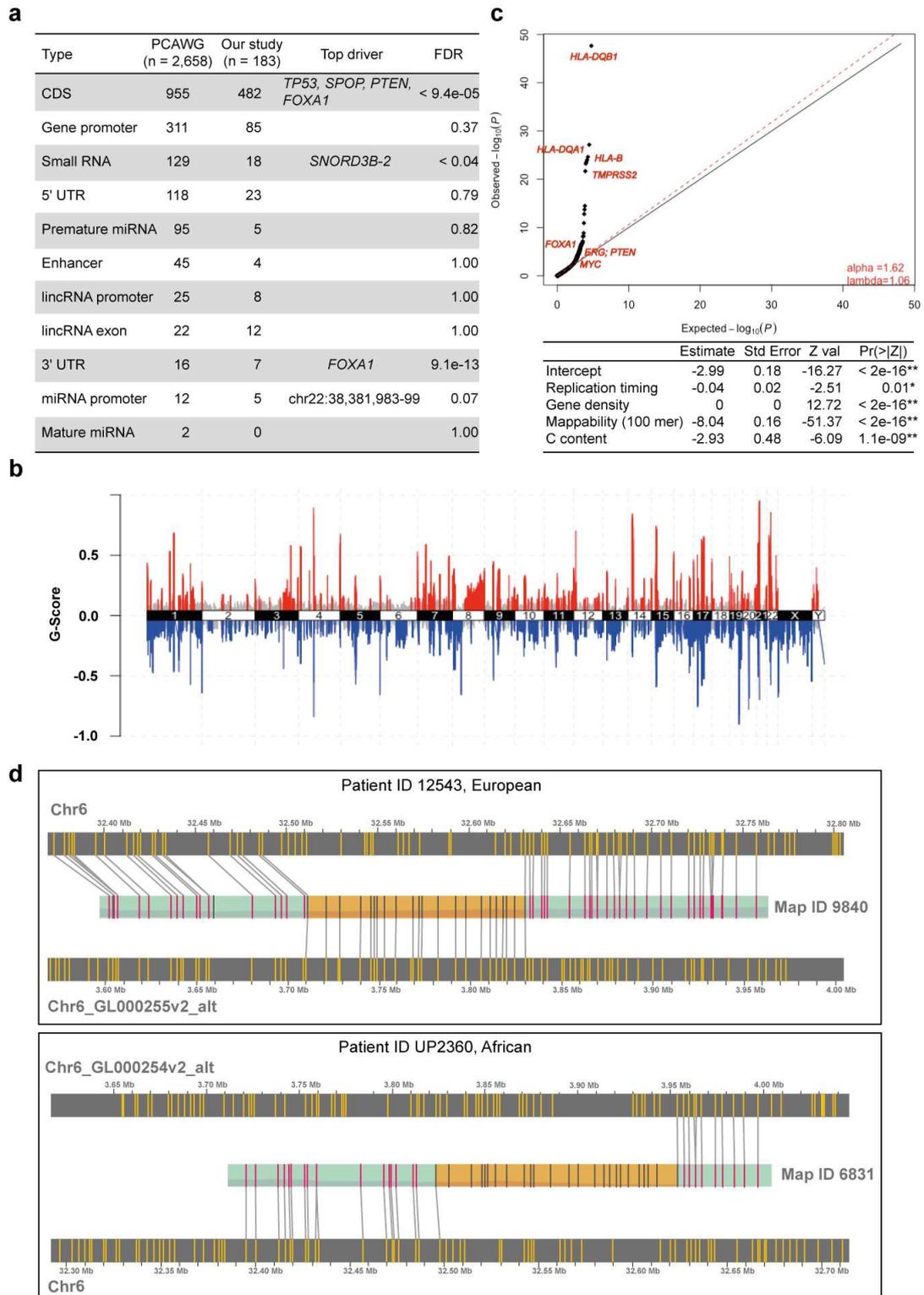
735

736 **Extended Data Fig. 1 | Clinical cohorts and statistical metrics.** **a**, Clinical and pathological patient
 737 characterisation. **b**, STRUCTURE analysis of bi-allelic germline variants with the logistic prior model.
 738 Model components used to explain structure in the plot are $K=5$. All spectrum of African contributions
 739 are summed and assigned as African ancestry. **c**, Saturation curve for all driver types across 183
 740 patients. Recurrent copy number gains and losses were measured using GISTIC v2 (Supplementary
 741 Methods). CDS, coding sequence; SV, structural variation. **d**, Spearman's correlation between different
 742 variables measured in this cohort. Dot sizes represent the magnitude of correlation, with significant P -
 743 values <0.01 .



744

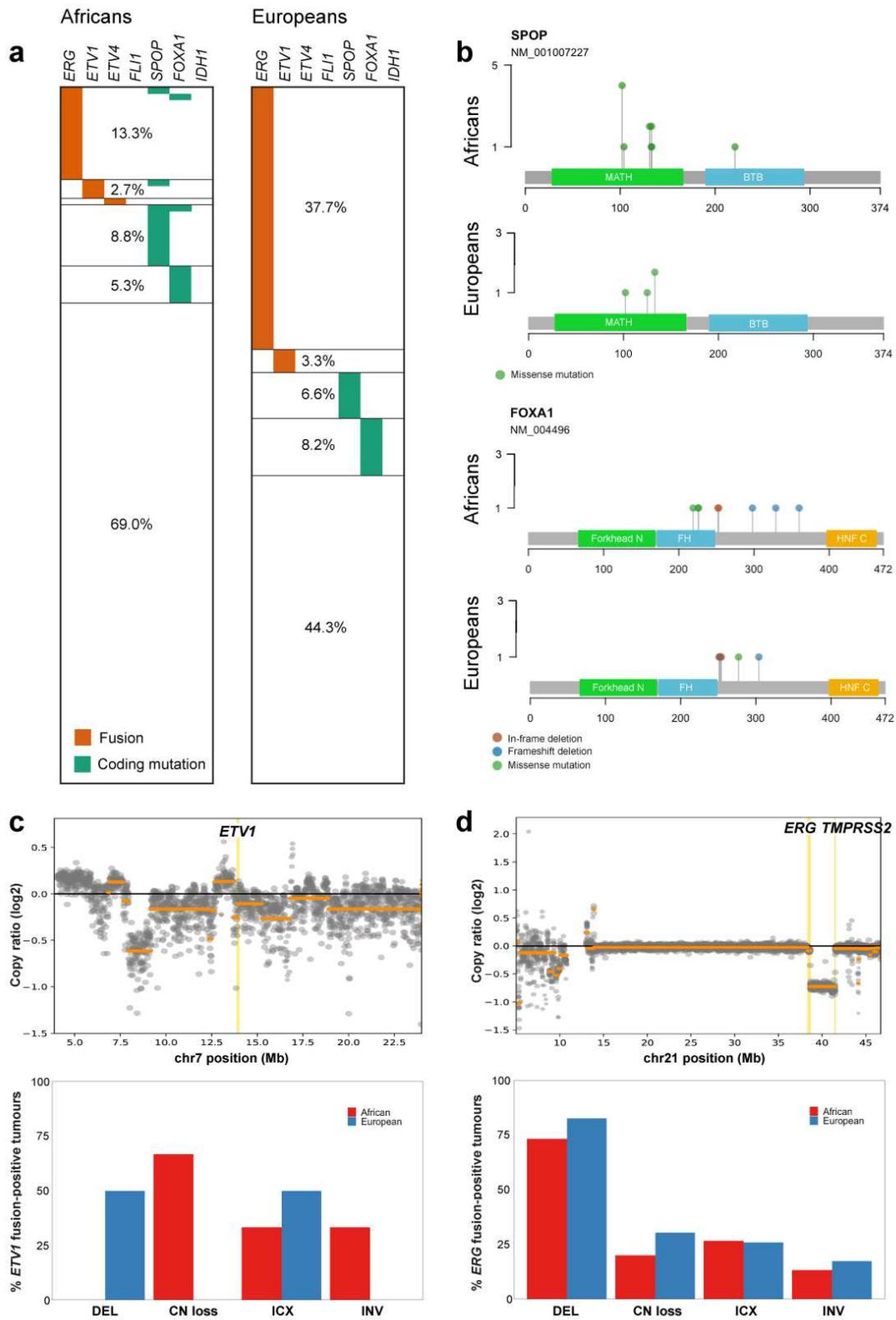
745 **Extended Data Fig. 2 | Somatic driver mutations in 183 prostate cancer patients.** The covariates
 746 on the left show mutational types and statistical significance (FDR) from ActiveDriverWGS and
 747 GISTIC2. **a**, The top 300 driver genes in PCAWG discovered in primary prostate tumours among 183
 748 specimens. The top barplot shows the distribution of the number of prostate cancer drivers and/or that
 749 of PCAWG. The heatmap shows drivers found in this study (rows) for each patient (columns).
 750 Heatmaps are coloured by mutational type. Bottom covariates show the clinical features of patients.
 751 The percentage of transition/transversion mutations across 183 patients shows 1,364,210 small somatic
 752 mutations across chromosomes 1-Y. **b**, The bottom heatmap shows the top 75 of previously reported
 753 coding driver genes in prostate cancer observed in this study^{7,8,18,19}. The right barplot shows the number
 754 of patients for each driver.



755

756 **Extended Data Fig. 3 | Discovery of prostate cancer drivers.** **a**, The number and types of PCAWG
 757 driver genes and elements studied in our cohort. **b**, Recurrent copy number alterations among 183
 758 prostate tumours identified with a 99% confidence level using GISTIC v2 (Supplementary Methods).

759 The figure shows GISTIC peaks of significant regions of recurrent amplification (red) or deletion
760 (blue) supported by FDR <0.01. **c**, Genome-wide scan for significantly recurrent breakpoints in our
761 study. The quantile-quantile plot shows *P*-values for mutational densities across 183 prostate cancer
762 patients. Generalised linear modelling (GLM) of somatic mutation densities along the genome with
763 significant background mutational processes adjusted in the model is also shown. **d**, Bionano
764 Genomics optical genome mapping at the HLA complex. Examples of HLA translocations from a
765 European patient (ID 12543) and an African patient (ID UP2360) studied in this cohort are
766 characterised by pairs of optical maps, each carrying a fusion junction with flanking fragments aligning
767 to one side of the two reference breakpoints. Using the recurrent HLA breakpoints identified in this
768 study, the genome map of the African specimen is found to have a low-end fusion function matched
769 with chromosome 6 through a manual inspection of unfiltered consensus maps using Bionano Access
770 v15.2. Note that the HLA alternate contig fused in the European tumour is different from one suggested
771 by short-read sequencing (chr6_GL000252v2_alt). The reference genome map is an *in silico* digest of
772 the human reference hg38 with the DLE-1 enzyme. Genome map sizes are indicated on the horizontal
773 axis, in megabase (Mb) units. Matching fluorescent labels between sample and reference genome map
774 are connected by gray lines.

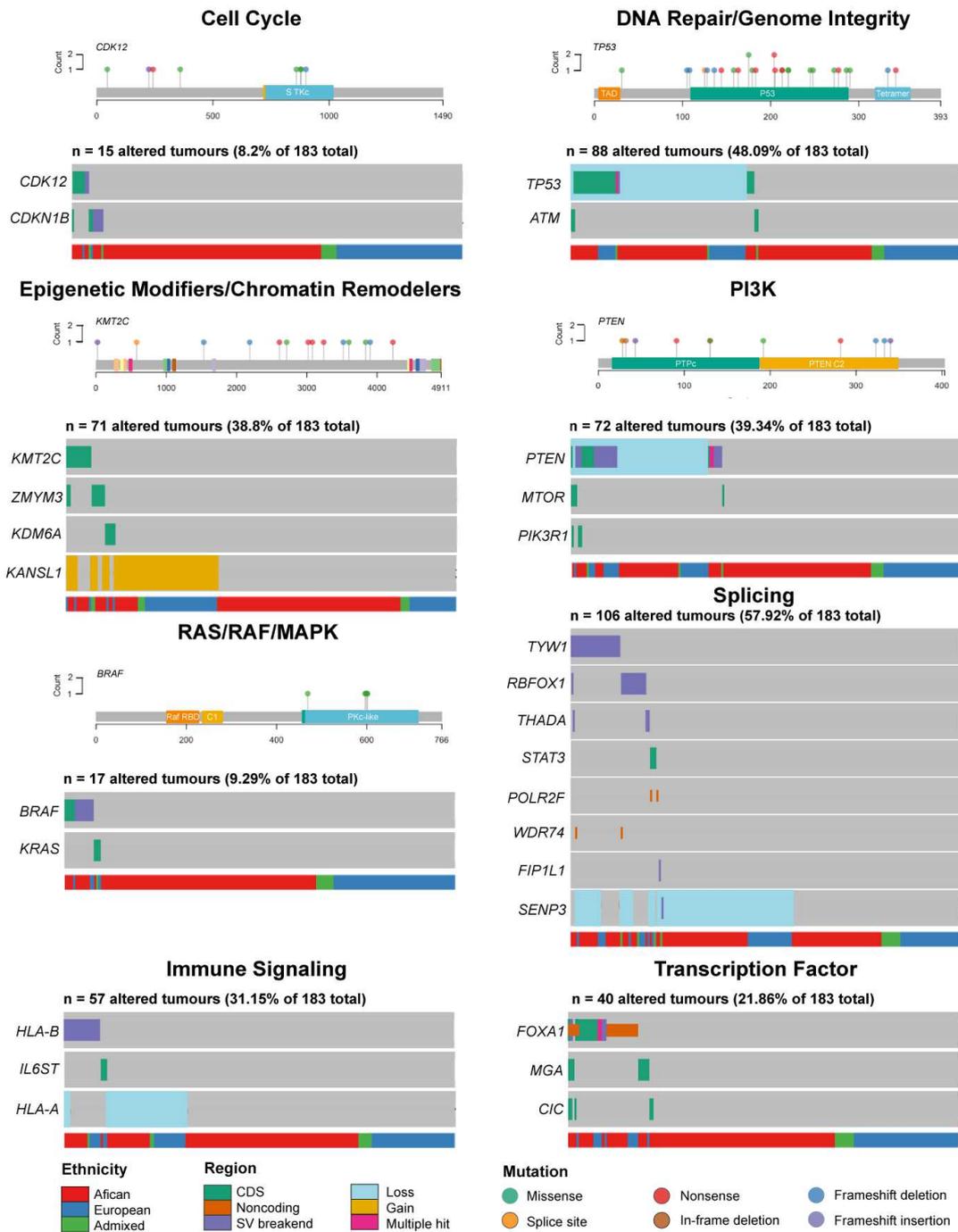


775

776 **Extended Data Fig. 4 | TCGA molecular taxonomy.** **a**, Seven important oncogenic drivers identified
 777 by TCGA within our African and European patients. **b**, Coding mutations observed within *SPOP* and
 778 *FOXA1* genes. Rarely, a mutation at the BTB domain of *SPOP* gene is shown (R221C in an African

779 patient, KAL0072). FH, forkhead. **c**, *ETV1* fusions within positive patients caused by copy number
780 (CN) losses and/or structural variants (DEL, deletion; ICX, interchromosomal translocation; and INV,
781 unbalanced or balanced inversion). CN changes in chromosome 7 show the *ETV1* loss with log₂ CN
782 ratio less than -0.2. **d**, *ERG* fusions caused by CN losses and/or structural variants.

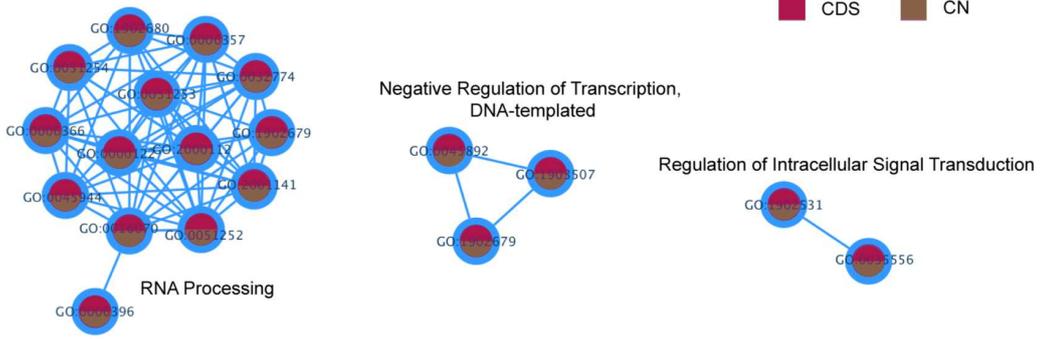
783



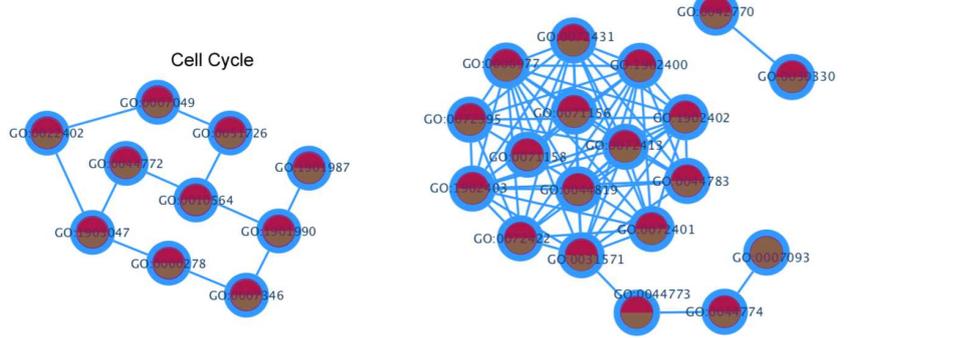
784

785 **Extended Data Fig. 5 | Prostate cancer genes and pathways.** The search is carried out using the
 786 TCGA and ICGC cancer databases. The top affected genes for each pathway are present with lollipop
 787 plots to show their hotspots of simple coding mutations if they existed.

a Signal Transduction



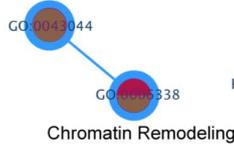
DNA Checkpoint Process



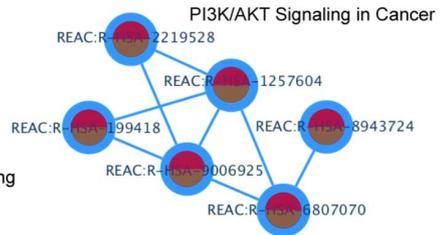
Immune Signaling



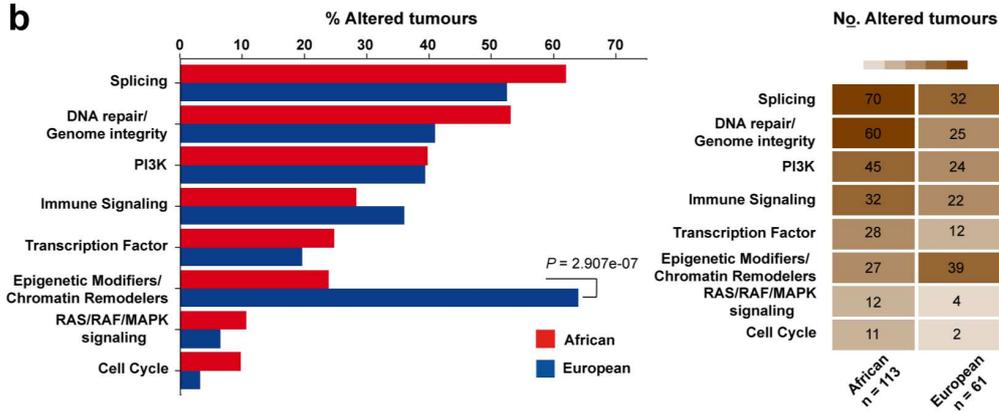
Chromatin Remodeling



PI3K/AKT Signaling



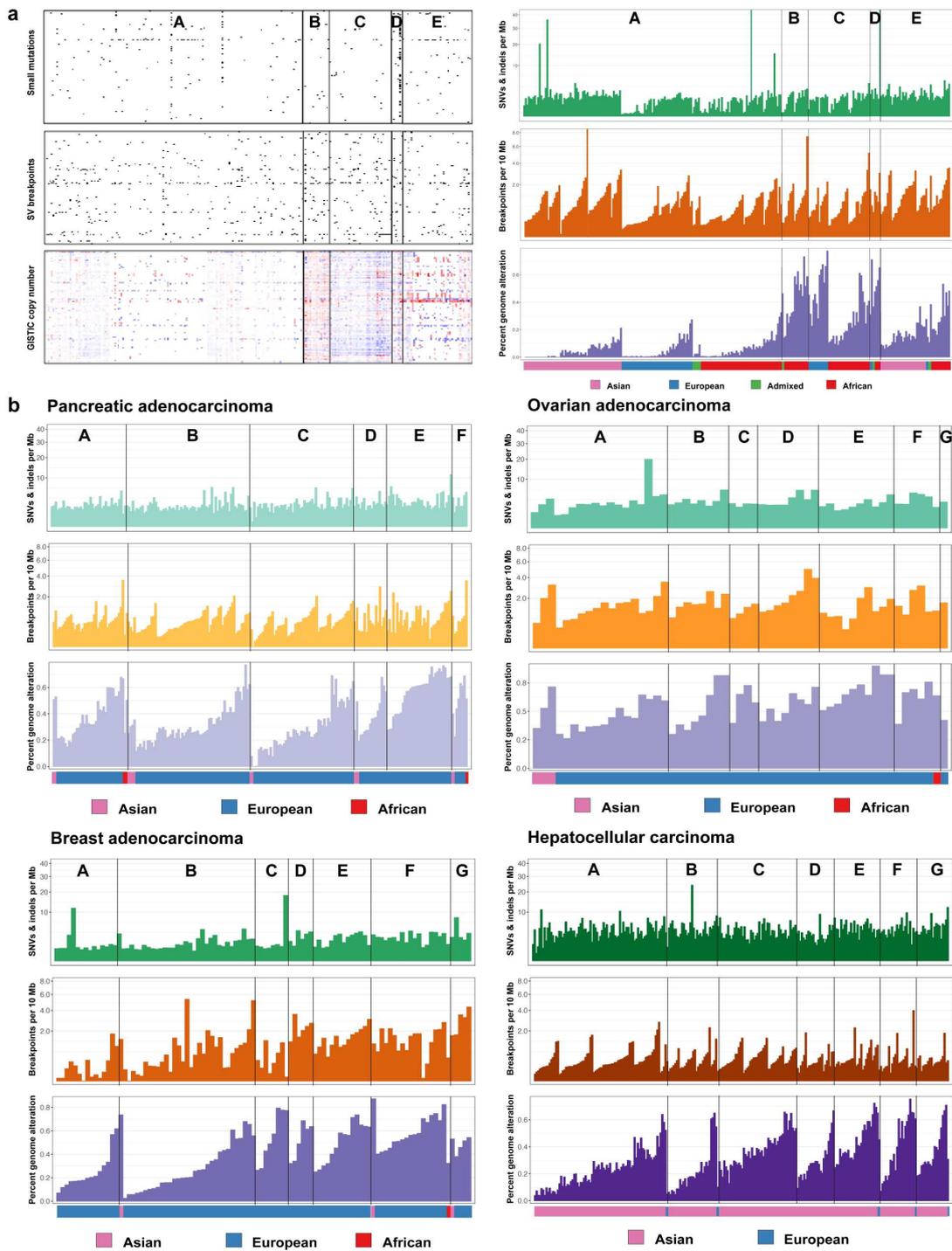
b



788

789 **Extended Data Fig. 6 | Major biological pathways and networks of prostate cancer. a,** Networks
 790 of functional interactions between driver genes are shown for each cancer pathway. Nodes represent
 791 Gene Ontology biological processes and Reactome pathways and edges show functional interactions.
 792 **b,** Pathway alteration frequencies between African and European. A sample was considered altered in

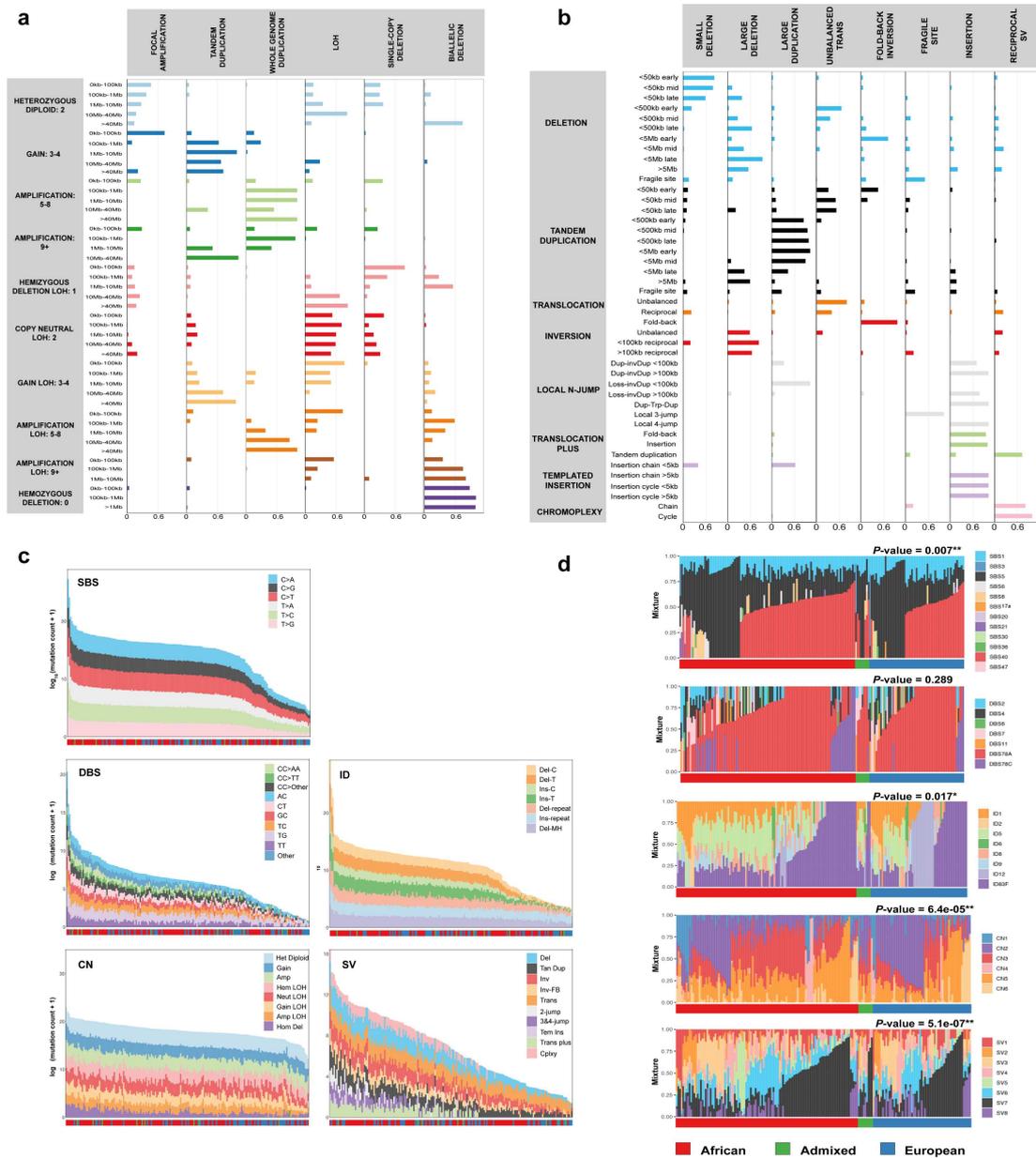
793 a given pathway if at least a single gene in the pathway had a genomic alteration. *P*-values indicate the
794 level of significance (two-sided Fisher's exact test).



795

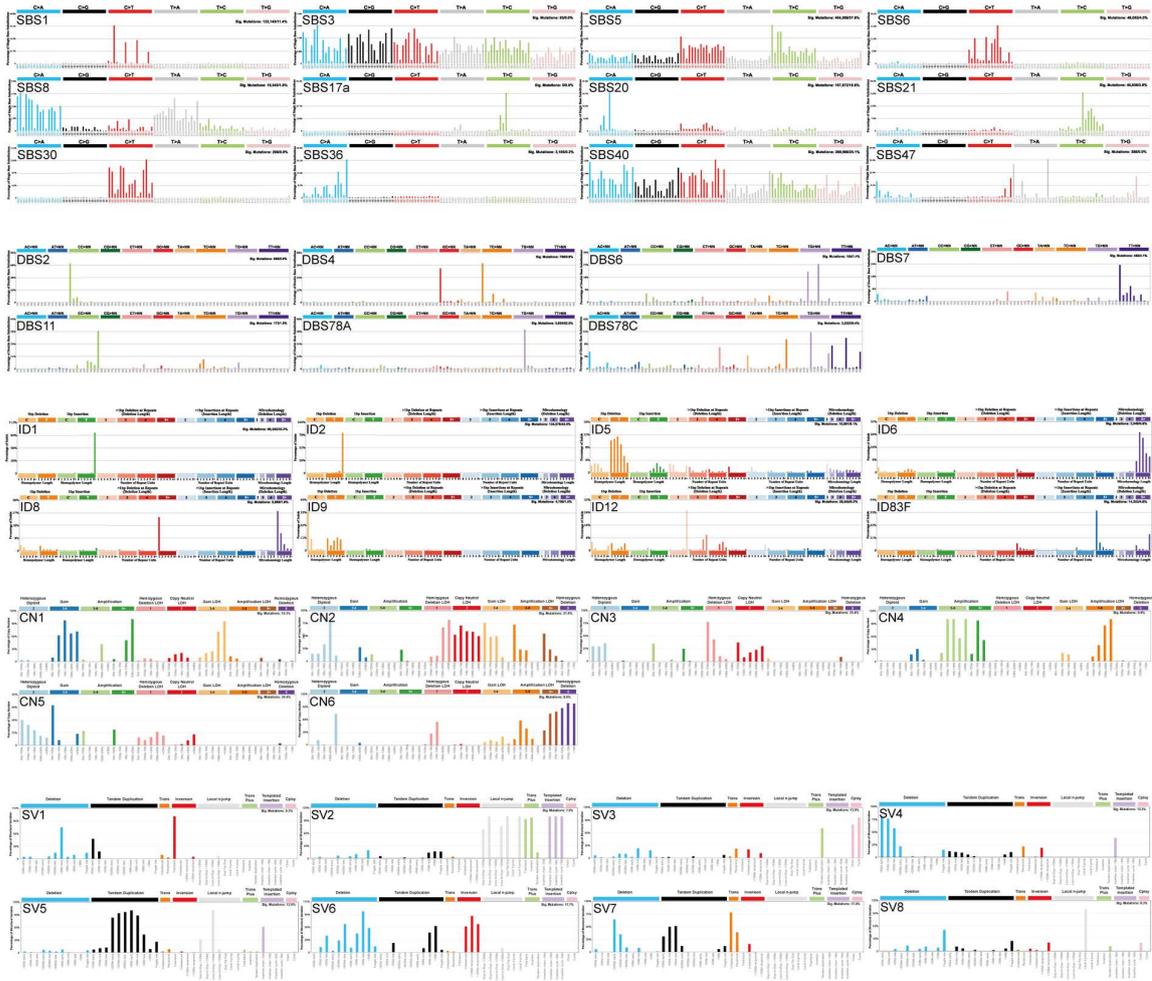
796 **Extended Data Fig. 7 | Molecular subtypes in prostate cancer and pan-cancers. a,** Unsupervised
 797 hierarchical clustering of primary prostate tumours across three major ethnic groups was performed
 798 using total somatic mutations present within WGS normalised data. Admixed individuals were also
 799 tested in prostate cancer subtypes to which they belonged. **b,** Molecular subtyping of total somatic
 800 mutations within pan-cancer studies, namely pancreatic, ovarian, breast and liver cancers. Raw data of

801 small somatic mutations, structural variants and copy number alterations acquired per cancer were
802 retrieved from the PCAWG¹⁴. For each subtype, patients are ordered based on their ethnicity. Ethnic
803 groups are assigned using a cut-off of ancestral contribution greater than 70%; otherwise, considered as
804 Admixed.



815 DBS, ID, CN and SV features across 183 tumours. Colours at the bottom panel show the following
816 ethnic groups: *i*) African, red; *ii*) Admixed, green; and *iii*) European, blue. **d**, Stacked barplots of
817 multiple signature exposures for each mutational type enriched per patient and ranked by ethnic group.
818 Copy number and structural variation signatures (CN1-6 and SV1-8, respectively) are the first
819 identified in this study for prostate cancer, and their enrichment in a patient appears to be significantly
820 associated (P -values <0.05) with our GMS, considering either *de novo* or global mutational signatures
821 discovered in the Catalogue of Somatic Mutations in Cancer (COSMIC).

822

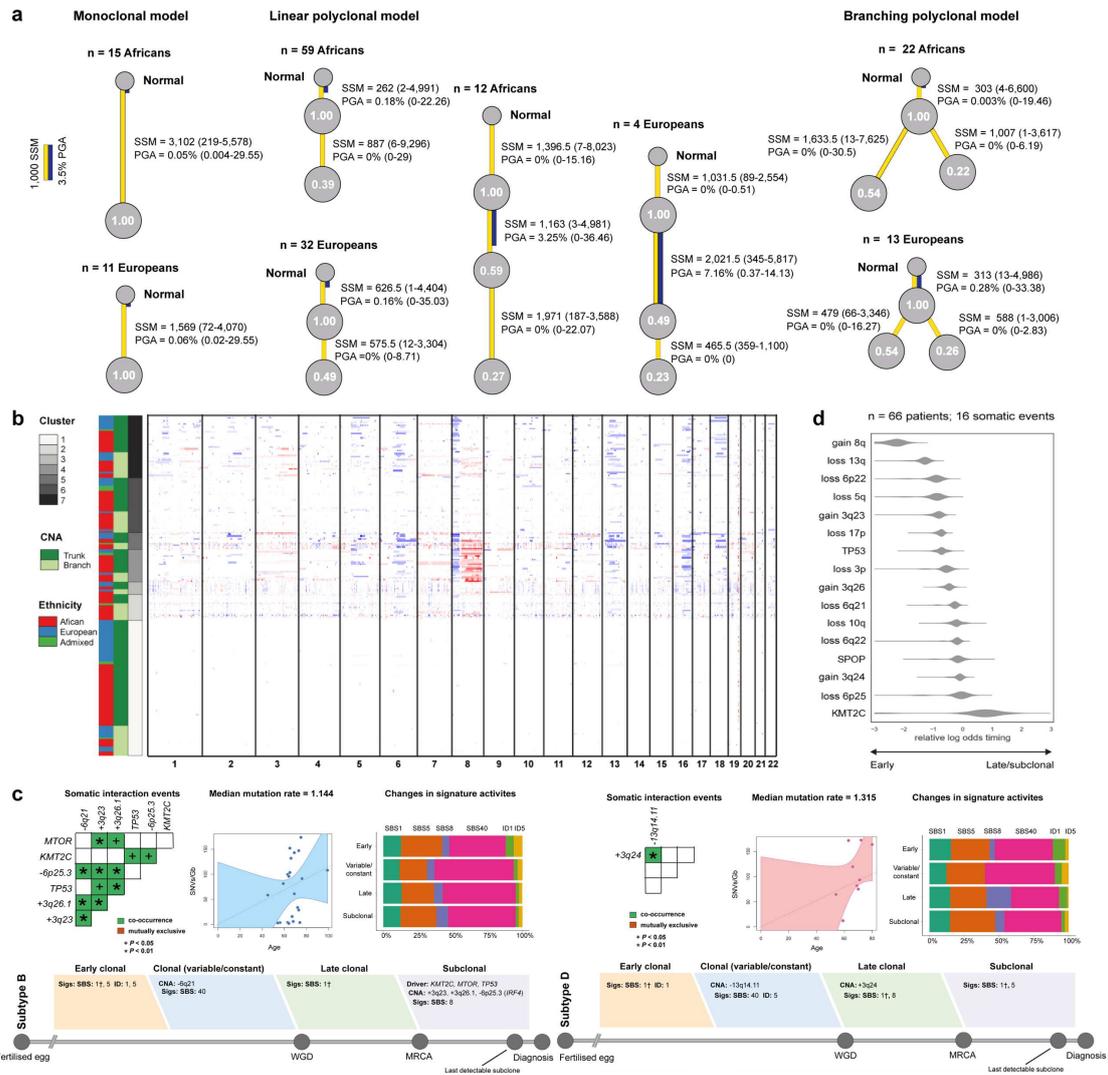


823

824 **Extended Data Fig. 9 | Total profiles of SBS, DBS, ID, CN and SV signatures.** The classification of
 825 each signature type (SBS, 96 classes; DBS, 78 classes; ID, 83 classes; CN, 45 classes; and SV, 44
 826 classes) is described in Supplementary Methods. The plotted data are available in digital form
 827 (Supplementary Table 9).

828

829



830

831 **Extended Data Fig. 10 | Stages of prostate tumour development. a**, Clonal architecture and its
 832 frequency in prostate cancer between Africans and Europeans. Tumours are divided into three groups:
 833 monoclonal, linear and branching polyclonal. The number of small somatic mutations (SSM) and CNA
 834 as percentage of genome alteration (PGA) is provided as median and range in bracket. Cancer cell
 835 fraction (CCF) in each clone and/or subclone is shown in a circular node. Tumours that show
 836 characteristics consistent with being polytumours or with multiple independent primary tumors are
 837 excluded to remain conservative. **b**, Unbiased hierarchical clustering of CNA between clonal (trunk)
 838 and subclonal (branch) mutations. Trunk mutations encompass those that occur between the root node
 839 (normal) and its only child node, while all others are classified to have occurred in branch. Red
 840 indicates gain; blue indicates loss; and rows indicate patients. Unidentified regions in trunk and branch
 841 are assumed to have neutral copy number. ConsensusClusterPlus showed seven CNA clusters among

842 our patients to be optimal. The figure shows that a trunk alteration from one patient is mutationally
843 similar to a branch alteration from another, rather than to other trunk ones from different patients in a
844 cohort. **c**, Cancer timelines of GMS-B and D identified in this study. Detailed explanation is provided
845 in Fig. 5. **d**, Relative ordering model (PhylogicNDT LeagueModel) results for a cohort of samples
846 (n=66). The samples can be analysed if they have somatic events of interest prevalent greater than 5%
847 of the sample size and have informative clonal status available for each event (16 events). Probability
848 distributions show the uncertainty of timing for specific events in the cohort.

849

850 **References**

- 851 1 Sung, H. *et al.* Global Cancer Statistics 2020: GLOBOCAN Estimates of
852 Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. *CA*
853 *Cancer J Clin* **71**, 209-249 (2021).
- 854 2 Alexandrov, L. *et al.* Signatures of mutational processes in human cancer.
855 *Nature* **500**, 415-421 (2013).
- 856 3 Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human
857 cancer. *Nature* **578**, 94-101 (2020).
- 858 4 Sandhu, S. *et al.* Prostate cancer. *Lancet* **398**, 1075-1090 (2021).
- 859 5 Boutros, P. C. *et al.* Spatial genomic heterogeneity within localized, multifocal
860 prostate cancer. *Nat Genet* **47**, 736-745 (2015).
- 861 6 Berger, M. F. *et al.* The genomic complexity of primary human prostate
862 cancer. *Nature* **470**, 214-220 (2011).
- 863 7 The-Cancer-Genome-Atlas-Network. The molecular taxonomy of primary
864 prostate cancer. *Cell* **163**, 1011-1025 (2015).
- 865 8 Wedge, D. C. *et al.* Sequencing of prostate cancers identifies new cancer
866 genes, routes of progression and drug targets. *Nat Genet* **50**, 682-692 (2018).
- 867 9 Lalonde, E. *et al.* Tumour genomic and microenvironmental heterogeneity for
868 integrated prediction of 5-year biochemical recurrence of prostate cancer: a
869 retrospective cohort study. *Lancet Oncol* **15**, 1521-1532 (2014).
- 870 10 Kamoun, A. *et al.* Comprehensive molecular classification of localized
871 prostate adenocarcinoma reveals a tumour subtype predictive of non-
872 aggressive disease. *Ann Oncol* **29**, 1814-1821 (2018).
- 873 11 Yamaguchi, T. N. *et al.* Molecular and evolutionary origins of prostate cancer
874 grade. .

- 875 12 Li, J. *et al.* A genomic and epigenomic atlas of prostate cancer in Asian
876 populations. *Nature* **580**, 93-99 (2020).
- 877 13 Crumbaker, M. *et al.* The Impact of Whole Genome Data on Therapeutic
878 Decision-Making in Metastatic Prostate Cancer: A Retrospective Analysis.
879 *Cancers (Basel)* **12**, E1178 (2020).
- 880 14 ICGC/TCGA-Pan-Cancer-Analysis-of-Whole-Genomes-Consortium. Pan-
881 cancer analysis of whole genomes. *Nature* **578**, 82-93 (2020).
- 882 15 Rotimi, S. O., Rotimi, O. A. & Salhia, B. A Review of Cancer Genetics and
883 Genomics Studies in Africa. *Front Oncol* **10**, 606400 (2021).
- 884 16 Jaratlerdsiri, W. *et al.* Whole Genome Sequencing Reveals Elevated Tumor
885 Mutational Burden and Initiating Driver Mutations in African Men with
886 Treatment-Naïve, High-Risk Prostate Cancer. *Can Res* **78**, 6736-6746 (2018).
- 887 17 Tindall, E. A. *et al.* Clinical presentation of prostate cancer in black South
888 Africans. *Prostate* **74**, 880-891 (2014).
- 889 18 Robinson, D. *et al.* Integrative clinical genomics of advanced prostate cancer.
890 *Cell* **161**, 1215-1228 (2015).
- 891 19 Armenia, J. *et al.* The long tail of oncogenic drivers in prostate cancer. *Nat*
892 *Genet* **50**, 645-651 (2018).
- 893 20 Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from
894 142 diverse populations. *Nature* **538**, 201-206 (2016).
- 895 21 Jaratlerdsiri, W. *et al.* KhoeSan Genome Project, a catalogue of ancient human
896 genome variation.
- 897 22 Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer
898 whole genomes. *Nature* **578**, 102-111 (2020).

- 899 23 Xia, L. *et al.* Multiplatform discovery and regulatory function analysis of
900 structural variations in non-small cell lung carcinoma. *Cell Rep* **36**, 109660
901 (2021).
- 902 24 Taylor, B. S. *et al.* Integrative genomic profiling of human prostate cancer.
903 *Cancer Cell* **18**, 11-22 (2010).
- 904 25 Gerstung, M. *et al.* The evolutionary history of 2,658 cancers. *Nature* **578**,
905 122-128 (2020).
- 906 26 Li, C. H., Haider, S. & Boutros, P. C. Ancestry Influences on the Molecular
907 Presentation of Tumours. *bioRxiv*.
- 908 27 Li, Y. *et al.* Patterns of somatic structural variation in human cancer genomes.
909 *Nature* **578**, 112-121 (2020).
- 910 28 Houlahan, K. E. *et al.* Germline determinants of the prostate tumor genome.
- 911 29 Schumacher, F. R. *et al.* Association analyses of more than 140,000 men
912 identify 63 new prostate cancer susceptibility loci. *Nat Genet* **50**, 928-936
913 (2018).
- 914 30 Al-Olama, A. A. *et al.* A meta-analysis of 87,040 individuals identifies 23 new
915 susceptibility loci for prostate cancer. *Nat Genet* **46**, 1103-1109 (2014).
- 916 31 Huang, F. W. *et al.* Exome Sequencing of African-American Prostate Cancer
917 Reveals Loss-of-Function ERF Mutations. *Cancer Discov*, doi:10.1158/2159-
918 8290 (2017).
- 919 32 Romanel, A. *et al.* Inherited determinants of early recurrent somatic mutations
920 in prostate cancer. *Nat Commun* **8**, 48 (2017).
- 921 33 Taylor, R. A. *et al.* Germline BRCA2 mutations drive prostate cancers with
922 distinct evolutionary trajectories. *Nat Commun* **8**, 13671 (2017).

- 923 34 Cairns, J. Mutation selection and the natural history of cancer. *Nature* **255**,
924 197-200 (1975).
- 925 35 Martincorena, I. & Campbell, P. J. Somatic mutation in cancer and normal
926 cells. *Science* **349**, 1483-1489 (2015).
- 927 36 Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic
928 cells. *Nat Genet* **47**, 1402-1407 (2015).
- 929 37 Ottman, R. Gene–Environment Interaction: Definitions and Study Designs.
930 *Prev Med* **25**, 764–770 (1996).
- 931 38 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-
932 Wheeler Transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 933 39 Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant
934 calls: the Genome Analysis Toolkit best practices pipeline. *Curr Protoc*
935 *Bioinformatics* **11**, 11.10.11-33 (2013).
- 936 40 Raj, A., Stephens, M. & Pritchard, J. K. fastSTRUCTURE: variational
937 inference of population structure in large SNP data sets. *Genetics* **197**, 573-
938 589 (2014).
- 939 41 Cortés-Ciriano, I. & Lee JJ, X. R., Jain D, Jung YL, Yang L, Gordenin D,
940 Klimczak LJ, Zhang CZ, Pellman DS; PCAWG Structural Variation Working
941 Group, Park PJ; PCAWG Consortium. Comprehensive analysis of
942 chromothripsis in 2,658 human cancers using whole-genome sequencing. *Nat*
943 *Genet* **52**, 331–341 (2020).
- 944 42 Baca, S. C. *et al.* Punctuated evolution of prostate cancer genomes. *Cell* **153**,
945 666-677 (2013).

- 946 43 Mermel, C. H. *et al.* GISTIC2.0 facilitates sensitive and confident localization
947 of the targets of focal somatic copy-number alteration in human cancers.
948 *Genome Biol* **12**, R41 (2011).
- 949 44 Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic
950 Tissues. *Cell* **171**, 1029-1041.e1021 (2017).
- 951 45 Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes
952 across 21 tumour types. *Nature* **505**, 495-501 (2014).
- 953 46 Mo, Q. *et al.* Pattern discovery and cancer gene identification in integrated
954 cancer genomic data. *Proc Natl Acad Sci U S A* **110**, 4245-4250 (2013).
- 955 47 Nik-Zainal, S. *et al.* Landscape of somatic mutations in 560 breast cancer
956 whole-genome sequences. *Nature* **534**, 47-54 (2016).
- 957 48 Du, Q. *et al.* Replication timing and epigenome remodelling are associated
958 with the nature of chromosomal rearrangements in cancer. *Nat Commun* **10**,
959 416 (2019).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [HRPCaSupplementaryMETHODS.pdf](#)
- [S1Clinicalcohortcharacteristicsandsequencingquality.xlsx](#)
- [S2Driverinformationbypatient.xlsx](#)
- [S3GISTIC2resultsofallgenomiclesionsunder99Xconfidencelevel.xlsx](#)
- [S4ListofsignificantlyrecurrentSVbreakpointsatFDRlowerthan0.10.xlsx](#)
- [S5TCGAprimatecancertaxonomyidentifiedinthisstudy.xlsx](#)
- [S6IntegrativeiClusteranalysisof183prostatetumours.xlsx](#)
- [S7Listof124preferentiallymutatedgeneswithinfourtumoursubtypes.xlsx](#)
- [S8Pathwayenrichmentanalysisof124preferentiallymutatedgenes.xlsx](#)
- [S9Totalmutationalsignatureprofilesacross183tumours.xlsx](#)
- [S10Crossindividualcontaminationlevel.xlsx](#)
- [S11Cancerevolutionanalysisofprostatecancer.xlsx](#)