

1 Sperm sequencing reveals extensive positive selection in the male germline

2 Matthew DC Neville¹, Andrew RJ Lawson¹, Rashesh Sanghvi¹, Federico Abascal¹, My H Pham¹, Alex
3 Cagan¹, Pantelis A Nicola¹, Tetyana Bayzetenova¹, Adrian Baez-Ortega¹, Kirsty Roberts¹, Stefanie V.
4 Lensing^{1,2}, Sara Widaa², Raul E Alcantara^{1,3}, María Paz García⁴, Sam Wadge⁴, Michael R Stratton¹,
5 Peter J Campbell¹, Kerrin Small⁴, Iñigo Martincorena¹, Matthew E Hurles¹, Raheleh Rahbari¹

6 ¹Cancer, Ageing and Somatic Mutation, Wellcome Sanger Institute, Hinxton, United Kingdom

7 ²Sequencing Operations, Wellcome Sanger Institute, Hinxton, United Kingdom

8 ³Quotient Therapeutics Limited, Saffron Walden, UK

9 ⁴Kings College London, Department of Twin Research & Genetic Epidemiology, London, United Kingdom

10

11

12

13 Abstract

14 Mutations that occur in the cell lineages of sperm or eggs can be transmitted to offspring. In humans,
15 positive selection of driver mutations during spermatogenesis is known to increase the birth prevalence
16 of certain developmental disorders. Until recently, characterising the extent of this selection in sperm
17 has been limited by the error rates of sequencing technologies. Using the duplex sequencing method
18 NanoSeq, we sequenced 81 bulk sperm samples from individuals aged 24 to 75 years. Our findings
19 revealed a linear accumulation of 1.67 (95% CI = 1.41-1.92) mutations per year per haploid genome,
20 driven by two mutational signatures associated with human ageing. Deep targeted and exome NanoSeq
21 of sperm samples identified over 35,000 germline coding mutations. We detected 40 genes (31 novel)
22 under significant positive selection in the male germline, implicating both activating and loss-of-
23 function mechanisms and diverse cellular pathways. Most positively selected genes are associated with
24 developmental or cancer predisposition disorders in children, while four genes that exhibit elevated
25 frequencies of protein-truncating variants in healthy populations. We find that positive selection during
26 spermatogenesis drives a 2-3 fold elevated risk of known disease-causing mutations in sperm, resulting
27 in 3-5% of sperm from middle-aged to elderly individuals carrying a pathogenic mutation across the
28 exome. These findings shed light on the dynamics of germline mutations and highlight a broader
29 increased disease risk for children born to fathers of advanced age than previously appreciated.

30

31

32 Introduction

33 Human cells in all tissues accumulate mutations throughout life. In replicating tissues, acquired driver
34 mutations that confer a selective advantage can promote the expansion of individual clones within
35 competing stem and progenitor cell populations. While patterns of selection and clonal expansion have
36 been extensively studied in cancers, recent research has also highlighted their occurrence in normal
37 tissues during ageing¹⁻¹⁰.

38

39 The spermatogonial stem cells of the testis occupy a unique niche amongst other studied normal tissues.
40 Among replicating cells, they have the lowest mutation rate, ~5-20 fold lower than any other studied
41 somatic cell type⁷. They are also the only replicating cells with the potential to transmit mutations to
42 offspring, balancing self-renewal and spermatogenesis to produce 150-275 million sperm per day post-
43 puberty^{11,12}. Targeted sequencing studies have revealed that driver mutations are acquired in
44 spermatogonial stem cells and that these cell populations expand along seminiferous tubules, resulting
45 in elevated fractions of mutant clones that are detectable in sperm¹³⁻¹⁷. Interestingly, all germline driver
46 mutations identified so far are activating missense hotspot mutations, which contrasts with a broader
47 range of activating and inactivating driver mutations observed in cancers and somatic tissues. These
48 germline driver mutations can have profound implications for offspring, as they are found in a set of 13
49 genes all known to cause severe developmental disorders¹⁸. This leads to a significant increase, up to
50 1,000-fold, in the sporadic birth prevalence of these disorders, with a strong correlation to elevated age
51 of the father¹⁹.

52

53 Technical limitations, related to the polyclonality and low mutation rate of testis and sperm, have
54 prevented extensive characterisation of this selection beyond a limited set of genes¹⁸. However, recent
55 advances in error-corrected duplex DNA sequencing approaches, in which information from both DNA
56 strands is used to detect mutations at single molecule resolution²⁰⁻²², have proven successful for the
57 accurate estimation of mutation burden in sperm²³⁻²⁵. Here we combine the duplex approaches of whole
58 genome NanoSeq²³ with deep whole exome and targeted NanoSeq (Lawson A.R., Abascal F., P.A.
59 Nicola et al., manuscript submitted for publication) to characterise positive selection in the male
60 germline and quantify its consequences for accumulation of disease mutations in sperm.

61

62 Results

63 Cohort and sequencing coverage

64 We performed restriction enzyme based, whole genome NanoSeq²³ of bulk semen samples (n = 81; 1-
65 2 timepoints per donor; age range: 24-75 years) and matched blood (n = 119; 1-3 timepoints; age range:
66 22-83 years) from 63 men in the TwinsUK cohort²⁶ (including 9 monozygotic and 3 dizygotic twin
67 pairs; **Methods; Supplementary Table 1**). The analysed sperm samples had sperm counts above 1

68 million/mL, as those below this threshold showed evidence of somatic cell contamination
69 (**Extended Data Fig. 1; Supplementary Note 1**). Across these samples, the mean number of unique
70 DNA molecules per site where a mutation was callable (duplex coverage - dx) was 3.7dx in sperm, and
71 4.3dx in blood (**Extended Data Fig. 2a**). For sperm, a haploid cell, 1dx is equivalent to one cell,
72 whereas for blood, a diploid cell, 2dx is equivalent to one cell.

73

74 **Mutational burden and signatures**

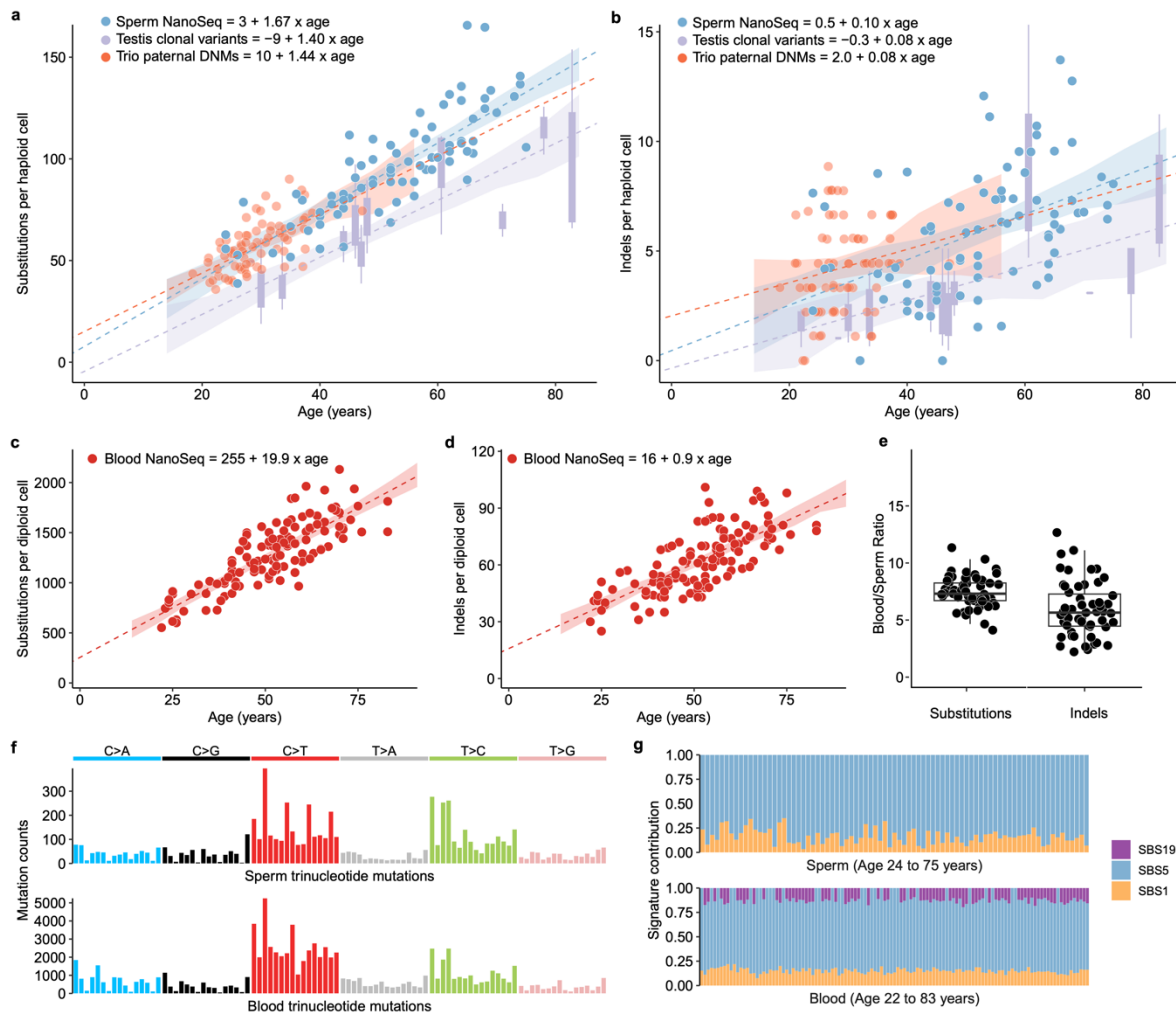
75 We performed stringent variant filtering of single nucleotide variants (SNVs) and small insertion–
76 deletion mutations (indels) from sperm and blood whole genome NanoSeq (**Methods**). From the 6,653
77 SNVs detected in sperm, we estimated an age-related accumulation of 1.67 substitutions per year per
78 haploid genome (95% CI 1.41-1.92, linear mixed-effect regression). This is comparable to estimates
79 from paternal *de novo* mutations (DNMs) in family pedigrees²⁷ of 1.44 substitutions per year (95% CI
80 1.00-1.87) and seminiferous tubules of the testes⁷ of 1.40 substitutions per year (95% CI 1.02-1.76;
81 **Fig. 1a**). Indels accumulated in sperm at a rate of 0.10 indels per year per haploid genome (95% CI 0.06-
82 0.15), similar to the rate observed in paternally phased DNMs²⁷ of 0.08 haploid indels per year (95%
83 CI -0.02-0.17) and seminiferous tubules of the testes⁷ of 0.08 haploid indels per year (95% CI 0.02-
84 0.13; **Fig. 1b**).

85

86 From the 92,035 SNVs and 4,641 indels detected in whole blood, we estimated an age-related
87 accumulation of 19.9 substitutions per year per diploid genome (95% CI 17.3-22.5; **Fig. 1c**) and 0.9
88 indels per year (95% CI 0.7-1.1; **Fig. 1d**). Both estimates are within the range of mutation rates observed
89 for specific cell types in the blood⁸, consistent with measuring a weighted average of these cell types in
90 whole blood (**Extended Data Fig. 3a,b**). We find that individuals had a mean of 7.6-fold more
91 substitutions per bp per year (range 4.2-11.5; **Fig. 1e**) and 6.3-fold more indels per bp per year (range
92 2.2-18.7; **Fig. 1e**) in blood than in sperm. Accounting for twin status or multiple timepoints from the
93 same individuals had a significant predictive effect for mutation burden in blood but not in sperm
94 (**Supplementary Note 2**).

95

96 The SNV mutational signatures in sperm were inferred to be SBS1 (mean 16%) and SBS5 (mean 84%),
97 the expected clock-like ageing signatures²⁸ (**Fig. 1f-g**). In blood, SBS1 (mean 15%) and SBS5 (mean
98 75%) were also the main mutational signatures, with an additional contribution of SBS19 (mean 10%),
99 which has been linked to persistent DNA lesions in hematopoietic stem cells²⁸ (**Fig. 1f-g**). We observed
100 that all signatures were correlated with age (**Extended Data Fig. 3c,d**). SBS1 and SBS5 accumulated
101 in individuals at a mean of 8.9-fold (range 2.3-39.1) and 6.8-fold (range 3.7-10.9) higher rate in blood
102 than in sperm respectively (**Extended Data Fig. 3e**), indicating that SBS19 does not explain a
103 substantial fraction of the mutation burden gap between the two tissues.



104

105 **Figure 1 | Mutational burden and signature analysis in sperm and matched blood**

106 **a,b**, Substitutions (**a**) and indels (**b**) per haploid cell from sperm whole genome NanoSeq, trio paternal DNMs²⁷
 107 called with standard sequencing and clonal variants from seminiferous tubules of testis⁷ called with standard
 108 sequencing. Dots indicate single donors while boxplots for testis variants show 1-15 samples per donor. **c,d**,
 109 Substitutions (**c**) and indels (**d**) per diploid cell for different ages from blood NanoSeq samples. **e**, Ratio blood to
 110 sperm substitutions and indels per diploid cell per year. Each dot corresponds to an individual with both a blood
 111 and sperm sample and where individuals had multiple timepoints the mean value of all timepoints in that tissue
 112 was used. **f**, Trinucleotide mutation counts in all sperm and blood samples. **g**, Contribution of signatures SBS1,
 113 SBS5, and SBS19 in sperm and blood samples ordered by age. **a,b,c,d**, Models are linear mixed regressions with
 114 95% CIs calculated by parametric bootstrapping. **a,b,e**, Box plots show the interquartile range, median, and 95%
 115 confidence interval for the median.

116

117 **Selective pressure dynamics in sperm**

118 To investigate positive selection in protein-coding regions in sperm we required much greater duplex
 119 coverage. Therefore, we utilised a capture-based modification to NanoSeq (Lawson A.R., Abascal F.,
 120 P.A. Nicola et al., manuscript submitted for publication) to deeply sequence coding regions from the

121 same set of semen samples. Specifically, we sequenced 38 samples using whole-exome NanoSeq to a
122 mean depth of 551dx per sample (20,923 cumulative dx), and 81 samples using targeted NanoSeq to a
123 mean depth of 985dx per sample (79,811 cumulative dx) with a target panel consisting of 263 canonical
124 cancer driver genes, 107 of which are also associated with developmental disorders (**Extended Data**
125 **Fig. 2a; Supplementary Table 2; Methods**). After variant filtering (**Methods**), we detected 56,503
126 (58% within coding regions) SNV/indel mutations from the exome panel and 5,059 (58% within coding
127 regions) from the targeted cancer panel. The age correlation of mutation burden for exome and targeted
128 sample sets were consistent with whole genome NanoSeq after correcting for the relative trinucleotide
129 composition of sequencing coverage (**Extended Data Fig. 2b**).

130

131 The vast majority of variants (99.5%) were detected only in a single duplex molecule of a sample.
132 Similarly, in the 23 samples with two timepoints (mean 12.1 year gap), 99.3% of the 5,143 variant calls
133 from the first timepoint were not called in the second timepoint. These results are consistent with sperm
134 being a highly polyclonal collection of cells derived from a large population of spermatogonial stem
135 cell progenitors in the testis.

136

137 The exome-wide strength of positive selection in sperm was quantified by estimating the rate of non-
138 synonymous (N) relative to selectively neutral synonymous (S) mutations (dN/dS ratio, where dN/dS =
139 1.0 indicates neutrality). We employed the *dNdScv* algorithm, which by default calculates dN/dS while
140 adjusting for trinucleotide context and several gene-level genomic covariates that influence mutation
141 rate²⁹. We modified this algorithm in three ways: first, we adjusted for duplex sequencing coverage per
142 base to correct for differential coverage within and between genes; second, we incorporated an
143 adjustment for CpG methylation levels in the testis due to its significant influence on mutation rates;
144 and third, we switched from trinucleotide to pentanucleotide context to better account for the effects of
145 extended contexts on germline mutation rates³⁰. These modifications refined exome-wide dN/dS ratios
146 by resolving specific mutation rate biases but had minor effects on gene-level dN/dS ratios
147 (**Extended Data Fig. 4, Supplementary Note 3**).

148

149 Using this model, we estimated the dN/dS ratio in the exome-sequenced samples to be 1.07 (95% CI
150 1.04-1.10). This ratio implies that 6.5% (95% CI 3.8%-9.1%) of the observed non-synonymous
151 substitutions in sperm conferred a clonal advantage during spermatogenesis in this cohort. Splitting the
152 cohort into thirds by age, we find that the exome-wide dN/dS ratio increased with age. The ratio was
153 1.01 (95% CI 0.93-1.09) in 26-42 year olds, 1.03 (95% CI 0.97-1.10) in 43-58 year olds, and 1.09 (95%
154 CI 1.06-1.13) in 59-74 year olds (**Fig. 2f**). This suggests that the dN/dS ratio increases over male
155 lifespan and that the cohort wide dN/dS ratios presented here in part reflect the age distribution of
156 samples (age range 26-74; mean 53 years).

157

158 We next compared the dN/dS ratios across gene sets related to spermatogenesis expression³¹ (**Fig. 2g**).
159 We find that the gene sets with the highest dN/dS ratios are those which are highly expressed during
160 spermatogenesis (1.25, 95% CI 1.13-1.38) and most specific to differentiated spermatogonial stem cells
161 (1.11, 95% CI 1.05-1.17). In contrast, the genes which are unexpressed in spermatogenesis (0.98, 95%
162 CI 0.88-1.11) and the genes most specific to elongating spermatids (1.01, 95% CI 0.94-1.08) show
163 dN/dS ratios close to neutrality. These results are consistent with the understanding that excess
164 nonsynonymous mutations observed in sperm confer a competitive advantage earlier in their cell
165 lineage, specifically in the spermatogonial stem cells of the testis¹⁵.

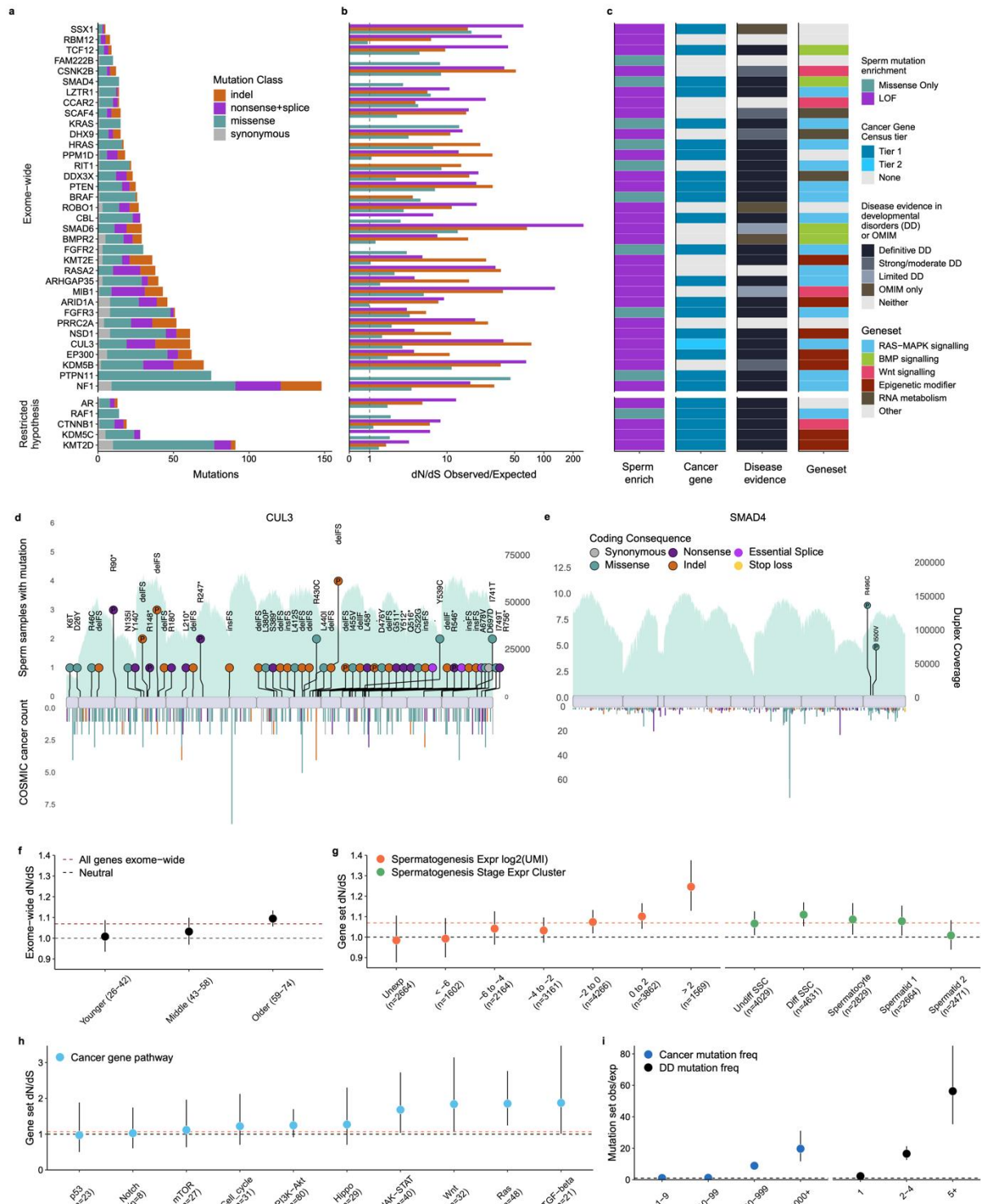
166

167 **Discovery of novel genes and pathways under positive selection in the germline**

168 We then investigated which genes were driving the signal of positive selection using the combination
169 of the exome and targeted panel datasets (**Methods**). We applied dN/dS tests for excess non-
170 synonymous mutations at both the gene-wide and SNV hotspot levels, which together identified 40
171 genes under significant positive selection. Of these, 35 genes reached exome wide significance at the
172 gene level (FDR $q < 0.1$; **Supplementary Table 3**) and/or contained one of 17 exome-wide significant
173 hotspots ($q < 0.1$; **Extended Data Table 1**). The genes *PTPN11*, *MIB1*, *RIT1*, *FGFR3*, *EP300*, and
174 *FGFR2*, were significant in both the gene and hotspot tests, *KDM5B*, *NF1*, *SMAD6*, *CUL3*, *RASA2*,
175 *PRRC2A*, *PTEN*, *ROBO1*, *DDX3X*, *CSNK2B*, *KRAS*, *PPM1D*, *ARID1A*, *BRAF*, *HRAS*, *KMT2E*,
176 *SCAF4*, *BMP2*, *TCF12*, *CCAR2*, *DHX9*, *NSD1*, *LZTR1*, *ARHGAP35*, *CBL*, *SSX1*, and *RBM12* were
177 significant in only the gene test, and *SMAD4* and *FAM222B* were significant in only the hotspot test
178 (**Fig. 2a,b**). We excluded the major seminal fluid component gene *SEMG1* from this list, despite excess
179 indels driving gene-level significance. This gene is expressed at extremely high levels in seminal
180 vesicles and unexpressed during spermatogenesis³⁶, suggesting that the enrichment may be the result of
181 a known process of indel hypermutation in highly-expressed genes^{37,38} from a small contamination of
182 seminal vesicle DNA, rather than selection in germ cells.

183

184 Subsequently, we carried out restricted hypothesis dN/dS tests at the per-gene and per-site level. The
185 gene level test examined only the set of 263 canonical cancer driver genes on our target panel and the
186 site level test used a set of 1,963 sites composed of known cancer hotspots and recurrent DNM sites
187 from the DDD cohort³⁹. This identified 5 additional genes: *KDM5C*, *KMT2D*, *AR*, *CTNNB1*, and *RAF1*
188 and 7 hotspots not already significant at the exome wide level ($q < 0.1$; **Fig. 2a,b**;
189 **Supplementary Table 4; Extended Data Table 1**).



190

191 **Figure 2 | Germline positive selection**

192 **a,b,c** Genes with significant dN/dS ratios from exome-wide and restricted hypothesis tests. **a**, Mutation count split
 193 by mutation class. **b**, Enrichment over expectation of mutation classes. **c**, Mutation type driving dN/dS enrichment,
 194 COSMIC³² cancer gene tier, developmental disorder gene link in DDG2P³³, and potential germline selection
 195 geneset. **d,e**, Observed sperm mutations across the cohort for *CUL3* and *SMAD4* where the height of the “lollipop”
 196 represents the number of unique samples with a mutation at that location and the colour represents its mutation

197 type. Mutations are labelled with their amino acid consequence for point substitutions or their insertion
198 (ins)/deletion (del) consequence of in frame (IF) or frameshift (FS). A “P” indicates that the variant is classified
199 as pathogenic/likely pathogenic in ClinVar³⁴. Exons are shown as purple rectangles and the blue background
200 represents the total duplex coverage across the cohort. Lines below the gene indicate COSMIC somatic mutations
201 in cancer within that gene³². **f-h**, dN/dS ratios for sperm SNVs across sets of individuals or genes, where the dotted
202 black line indicates neutrality and the dotted orange line represents the cohort average across all genes. **f**, Exome-
203 wide dN/dS ratios in sperm for the cohort split into thirds by age. **g**, Expression gene sets from single-cell
204 sequencing of germ cells³¹. Expression levels represent 7 bins of mean expression levels across germ cell stages
205 and expression clusters represent genes most characteristic to certain germ cell stages. Germ cell types include
206 undifferentiated and differentiated spermatogonial stem cells (SSCs), spermatocytes, round spermatids (1) and
207 elongating spermatids (2). **h**, Germline selection genes and cancer gene census genes split by ten canonical cancer
208 pathways in KEGG³⁵. **i**, The observed/expected mutation rate in sperm for bins of mutations. COSMIC and DDD
209 are bins of variants that have been seen at different levels of recurrence. Error bars depict 95% CIs.

210

211

212 Genes linked to germline positive selection to date all operate through activating missense mutations,
213 with 12 linked to the RAS-MAPK signalling pathway¹⁸ and one (*SMAD4*) linked to TGF- β /BMP
214 signalling⁴⁰. Our findings replicate *SMAD4* and 8 of the 12 RAS-MAPK pathway genes as under
215 significant positive selection in this dataset. The 4 genes which did not reach significance (*MAP2K1*,
216 *MAP2K2*, *SOS1*, and *RET*) each had between 2- and 4-fold enrichment of missense mutations, which
217 corresponded to nominally significant missense enrichment in all 4 genes ($p < 0.1$). Given the direct
218 evidence for these genes driving clonal selection in testis and nominal enrichment from sperm
219 sequencing, we expect that each will reach exome-wide significance with deeper sequencing.

220

221 We estimate that together, the 44 genes linked to germline selection here or in previous studies, contain
222 an estimated 357 (95% CI: 319–387) excess non-synonymous variants in exome sequenced samples.
223 This would account for 23% (95% CI: 14%–43%) of the total estimated driver variants across the
224 exome. The wide confidence intervals and the sensitivity of this estimate to the mutation model used
225 (**Supplementary Note 3**) suggest that small uncertainties in mutation rates, when propagated across the
226 exome, make it difficult to precisely estimate the fraction of drivers explained. Nevertheless, the
227 findings suggest that additional driver genes remain to be discovered.

228

229 The 31 newly identified genes demonstrate that germline positive selection is not restricted to activating
230 mutations or to the RAS-MAPK pathway. For instance, 30/31 of the novel genes are enriched for loss-
231 function mutations such as nonsense, splice, and indel variants, suggestive of protein-inactivating
232 mechanisms of selection (**Fig. 2c,d**). Splitting the germline selection genes and known cancer genes by
233 ten canonical cancer pathways within the Kyoto Encyclopedia of Genes and Genomes (KEGG)³⁵, we
234 find that the top gene groups enriched in dN/dS are RAS-MAPK, Wnt and TGF- β /BMP signalling

235 **(Fig. 2h)**. Indeed, many of the novel selection genes are linked to the RAS-MAPK pathway, such as
236 *NF1*, *CUL3* and *LZTR1*, Wnt signalling (i.e. *CSNK2B*, *MIB1*, *CCAR2*), and TGF- β /BMP signalling (i.e.
237 *SMAD6*, *TCF12*, *BMPR2*) **(Fig. 2c)**. We also identified a number of genes that encode epigenetic
238 modifiers (i.e. *KDM5B*, *KDM5C*, *ARID1A*, *KMT2D*, *KMT2E*, *EP300*, *NSD1*), and genes encoding RNA
239 metabolism proteins (i.e. *DHX9*, *DDX3X*, *SCAF4*). These findings highlight a new diversity of genes,
240 mutational mechanisms, and pathways driving germline selection, however future work will be needed
241 to confirm the specific pathways and roles through which these genes drive clonal expansion during
242 spermatogenesis.

243

244 It has been observed that cancers and germline developmental disorders share causal pathways and
245 genes⁴¹⁻⁴⁴. Notably, the 13 genes previously linked to germline positive selection are all known cancer
246 and known developmental disorders genes^{18,40}. This pattern holds, but to a lesser extent, in the new
247 germline selection genes identified here: 16 of 31 genes are tier 1 or 2 cancer census genes³² and 24 of
248 31 are linked to monogenic developmental disorders in the DDG2P database³³ **(Fig. 2c;**
249 **Supplementary Table 5)**. Among the positively selected genes which are not associated with
250 monogenic developmental disorders, three are linked to other monogenic disorders⁴⁵ **(Fig. 2c;**
251 **Supplementary Table 5)**. These include *BMPR2* associated with pulmonary arterial hypertension⁴⁶⁻⁴⁸,
252 *ROBO1* associated with pituitary stalk interruption syndrome^{49,50} and *SSX1* associated with X-linked
253 spermatogenic failure⁵¹.

254

255 The overlap between germline positive selection genes and cancer/developmental disorders is also
256 apparent at the variant level. Somatic mutations that are most frequently observed (>50 times) in the
257 COSMIC cancer database are enriched 11-fold (95% CI: 6-20) among our sperm mutation dataset after
258 adjusting for expected mutation rate **(Methods)**. Similarly, germline mutations that are most frequently
259 observed (>5 times) in a large cohort of children with developmental disorders are enriched 66-fold
260 (95% CI: 41-100) in our sperm mutation dataset **(Fig. 2i)**. In addition, the mutation types (e.g. missense
261 or protein truncating variants) enriched in sperm for a given gene are largely consistent with those
262 enriched in cancer and those causal for developmental disorders **(Extended Data Fig. 5a-c)**. These
263 results show a clear overlap between genes, hotspots, and mutation mechanisms which drive germline
264 positive selection, cancer, and developmental disorders. A notable exception to this pattern is *SMAD4*,
265 which has two distinct missense hotspots in sperm that are developmental disorder hotspots causal for
266 Myhre syndrome⁵² but that are not often seen in cancers, replicating recent findings⁴⁰ **(Fig. 2e)**.

267

268 **Positive selection drives enrichment of disease-causing mutations in sperm**

269 Given the association of many positively selected genes to disease, it is of interest to assess to what
270 degree germline positive selection may increase the fraction of sperm carrying potential disease-causing

271 mutations, and thus the birth prevalence of the associated disease. To estimate the fraction of sperm
272 carrying specific classes of variants, we aggregated the variant allele frequencies (VAFs) of different
273 mutation types and compared this to expected values. Expected values were generated using the custom
274 fit *dNdScv* mutation model (adjusted for base pair coverage, CpG methylation, and trinucleotide
275 context) and normalised to account for the linear impact of age on mutation rate (**Methods**).

276

277 We find that the fraction of sperm carrying a non-coding or synonymous mutation increases linearly
278 with age as predicted by the model (**Extended Data Fig. 6**). In contrast, the frequency of missense,
279 truncating (nonsense and splice) variants, and coding indels deviate above expected values in older
280 individuals, consistent with dN/dS results and indicative of positive selection acting over time
281 (**Extended Data Fig. 6**).

282

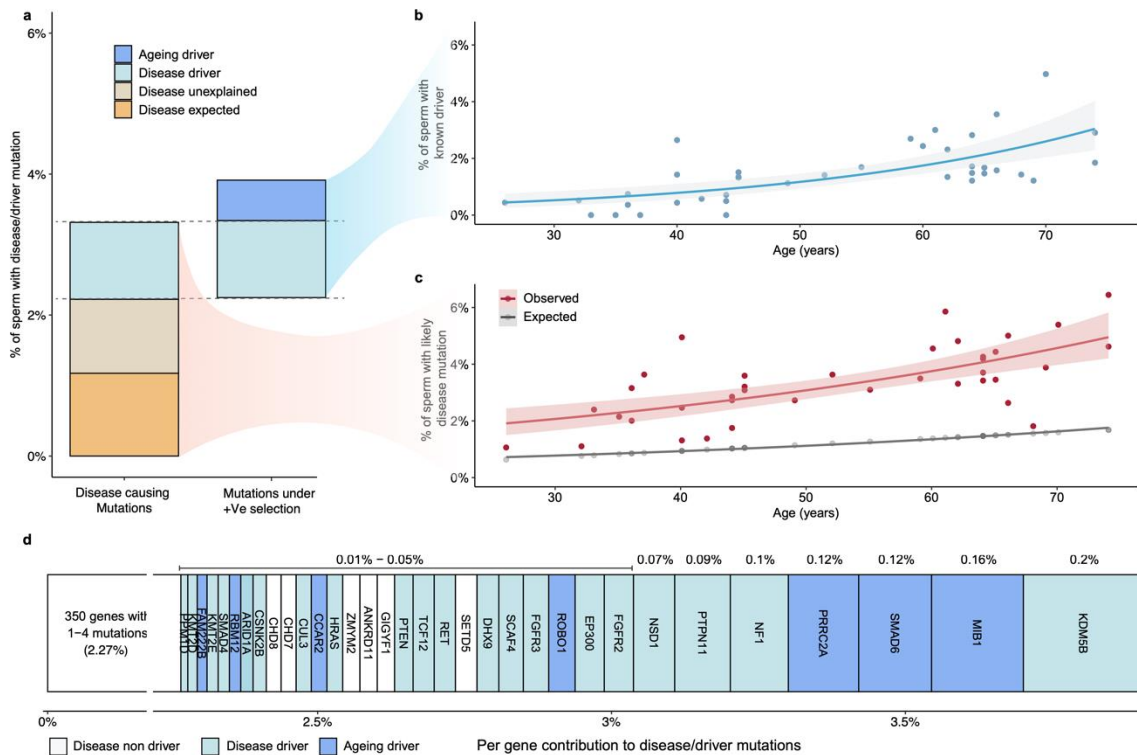
283 We then generated a list of likely monoallelic disease-causing mutations, which includes ClinVar³⁴
284 pathogenic/likely pathogenic variants and highly damaging variants in high confidence monoallelic
285 developmental disorder genes from DDG2P³³ (loss-of-function or missense CADD⁵³ score >30). This
286 list represents a conservative estimate of disease-causing mutations due to the incomplete discovery and
287 annotation of disease-causing variants and genes. We found that the observed fraction of sperm
288 containing disease mutations was markedly higher than that expected under a germline mutational
289 model at all ages of our cohort. The expected fraction of sperm with a likely disease mutation ranged
290 from 0.73% in 30-year-olds to 1.6% in 70 year olds, whereas, fitting a quasibinomial regression, the
291 observed fraction of sperm with a likely disease mutation in each age bracket ranged from 2% (95% CI:
292 1.6%-2.5%) in 30-year-olds to 4.5% (95% CI: 3.9%-5.2%) in 70 year olds (**Fig. 3c**). These differences
293 represent similar enrichments of 2.8-fold (95% CI: 2.2 to 3.5) and 2.9-fold (95% CI: 2.5 to 3.3) in 30-
294 year-olds and 70-year-olds respectively.

295

296 Interestingly, the disease cell fraction estimates are made up of many low frequency variants rather than
297 being driven by individual high VAF mutations. The estimates in exome samples are made up of a mean
298 of 18.3 unique variants (range 4-62) per individual. Furthermore, 692 of 696 (99.4%) of all those
299 variants are only observed in a single sperm cell, similar to the average of all variants (99.5%).

300

301 We next investigated to what degree the observed enrichment of disease mutations can be attributed to
302 driver mutations in positively selected genes. Fitting a quasibinomial regression, we observe a strong
303 positive correlation between age and driver rate ($P = 7.95e-06$) with an estimated 0.5% (95% CI: 0.3%-
304 0.8%) of sperm from individuals at age 30 and 2.6% (95% CI: 2.0%-3.3%) of sperm from individuals
305 aged 70 carrying a known driver mutation (**Fig. 3b**). However only about two thirds (65.6%) of those
306 driver mutations meet our criteria of likely disease-causing.



307

308 **Figure 3 | Pathogenic burden**

309 **a**, Estimated mean percentage of sperm in the cohort carrying a likely monoallelic disease mutation (Left) or a
 310 driver mutation in a germline selection gene (Right). Disease mutations are divided into the fraction that was
 311 expected from the mutation model, the portion explained by driver variants and the portion unexplained. Driver
 312 mutations are split by those contributing to the disease mutations and the remainder, 'Ageing drivers'. **b**,
 313 Estimated percentage of sperm per individual carrying a driver mutation by age. **c**, Observed and expected
 314 percentages of sperm with likely disease mutation by age. **b,c**, Model fits represent quasibinomial regressions
 315 with 95% confidence intervals. **d**, Cohort means from (**a**) split by gene and ordered by estimated mutation
 316 percentage. Per-gene contributions are shown above each gene; the summed contributions of all genes
 317 shown below. Genes with 4 or fewer variants are grouped on the left with a condensed x-axis for clarity.

318

319 Mutations in sperm that are likely disease-causing and those that are known driver mutations therefore
 320 represent overlapping but distinct annotations (**Fig. 3a**). Across the cohort, an estimated 3.3% of sperm
 321 carry a likely disease-causing mutation. Of this, approximately one-third (1.2%) is expected by the
 322 neutral mutation model, another third (1.1%) is explained by known driver mutations, and the remaining
 323 third (1.0%) is unexplained by either source. These findings suggest that the increase in likely disease-
 324 causing mutations is largely driven by germline positive selection, but also indicate that additional
 325 driver genes with disease associations remain to be identified.

326

327 Examining driver mutations which do not meet our criteria of a likely disease-causing mutation, we
 328 find that they impact an estimated 0.6% of sperm across the cohort. The consequences of these
 329 mutations are unclear. For instance, driver variants in *SMAD6*, which has a variably penetrant link to

330 congenital phenotypes⁵⁴, may be disease-causing in some cases but not others. Other potential
331 consequences include mutations that are disease-causing but not yet annotated as such, less able to
332 fertilise an egg, embryonic lethal, or biallelic disease-causing.

333

334 Much of the fraction of sperm with a disease and/or a driver mutation can be attributed to a small number
335 of genes in the exome. From 374 genes with at least one such variant, the 33 genes with ≥ 5 independent
336 mutations, most of which are under significant positive selection (26/33), represent 42.8% of the
337 disease/selection fraction in sperm (**Fig. 3d**). Strikingly, 6 of those genes, all of which are under
338 significant positive selection, (*KDM5B*, *MIB1*, *SMAD6*, *PRRC2A*, *NFI*, and *PTPN11*) together explain
339 over 20% of the disease/selection fraction. This suggests that although individual mutations we observe
340 are at low frequencies, positive selection systematically favours the likelihood of observing variants in
341 driver genes.

342

343 We next sought to examine whether there are risk factors other than age which contribute to the
344 accumulation disease or driver mutations in sperm. Currently known mutagenic effects in the male
345 germline include chemotherapy and inherited DNA repair defects^{39,55} and small effect size influences
346 of genetic ancestry and smoking⁵⁶. While the cohort did not include any individuals with known
347 chemotherapy treatments or DNA repair defects, phenotype data was available for BMI, smoking, and
348 alcohol consumption, all of which have some evidence of driving mutation burden or driver mutation
349 rate in some somatic tissues⁵⁷. We used multivariate generalised linear models to test the association
350 between these factors and measures of mutation burden, signatures, and driver cell fractions, correcting
351 for multiple testing (**Methods; Extended Data Fig. 6**). Regardless of the sperm sequencing set
352 examined (targeted, exome, or whole genome) only age was significantly correlated with measures of
353 mutational burden, signatures, and driver cell fractions. No significant effects were found for BMI,
354 smoking pack-years, or alcohol drink-years. However, in blood samples, age, smoking, and alcohol
355 consumption showed significant effects on SNV and SBS5 burdens. These results suggest that, unlike
356 many somatic tissues, the male germline mutation landscape may be largely protected from these risk
357 factors, although larger cohorts will be needed to interrogate possible small effect sizes.

358

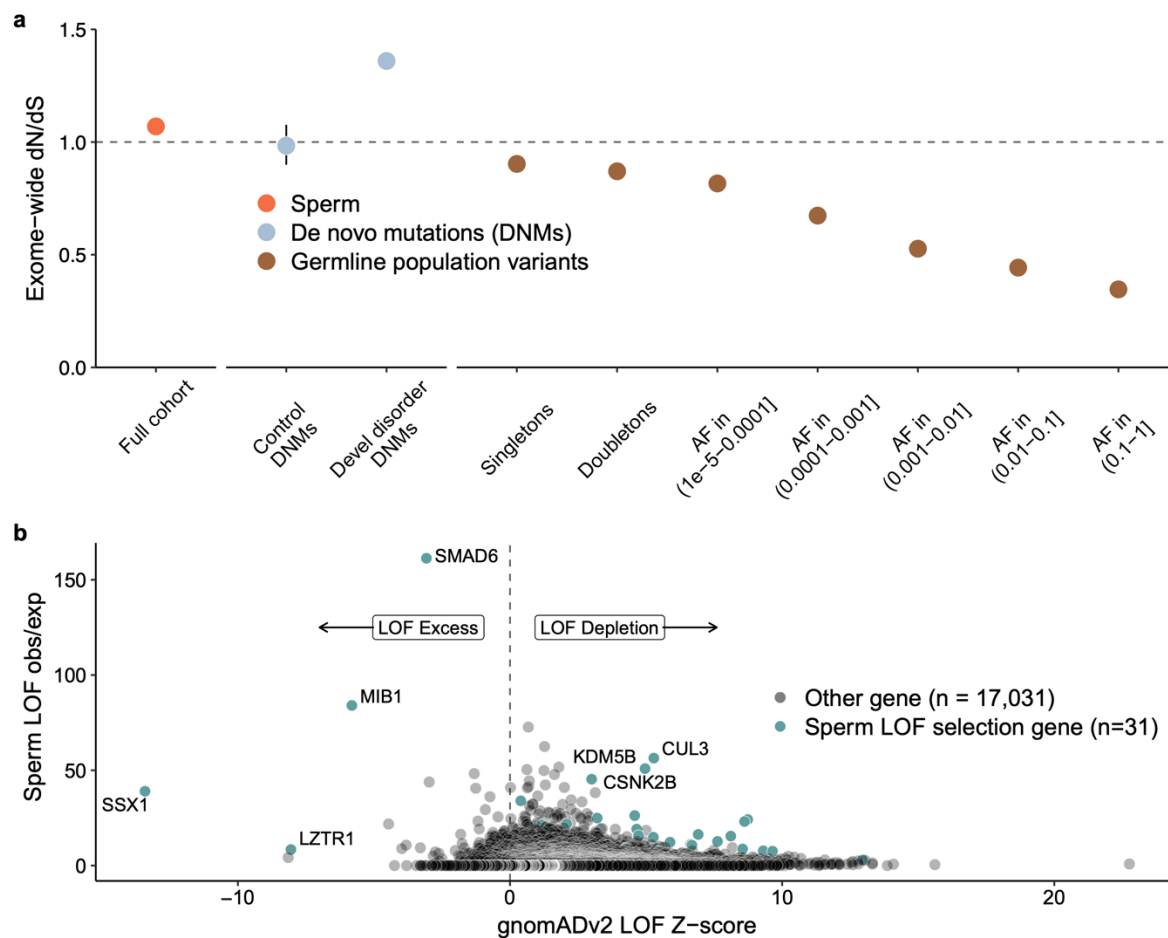
359 **Selection in germ cells relative to single generation and population-level variants**

360 Mutations in sperm account for ~80% of DNMs and are therefore also the origin of most population
361 level variants. Comparisons between these different sources of germline variants provide an opportunity
362 to explore how positive selection in the male germline may manifest over time.

363

364 Examining control DNMs from offspring without clinical phenotypes⁵⁸ we found a dN/dS ratio
365 consistent with neutrality of 0.98 (95% CI 0.90-1.08; **Fig. 4a**). In contrast, the dN/dS ratio in DNMs

366 from offspring with developmental disorders showed a large enrichment of nonsynonymous variants,
 367 as previously reported³⁹: 1.36 (95% CI 1.33-1.39; **Fig. 4a**). However, ascertainment biases in these
 368 cohorts suggest that these dN/dS ratios may not accurately reflect the dN/dS ratio of DNMs entering
 369 the population. Future large DNM studies of birth cohorts will likely be required to give less biased
 370 estimates.



371
 372 **Figure 4 | Comparison to population variation**

373 **a**, Exome-wide dN/dS ratios across different variant sets, including sperm variants from all exome sequenced
 374 samples, de novo mutations (DNMs) from a collection of healthy trios⁵⁸ and the Deciphering Developmental
 375 Disorders (DDD) cohort³⁹; and population variants from gnomAD⁵⁹ split by allele frequency (AF). **b**,
 376 Observed/expected enrichment of missense and loss-of-function (LOF; essential splice, nonsense, indel) variants
 377 in positively selected genes within sperm (x-axis) from dN/dS models vs gnomAD LOF z-score. Positive z-
 378 scores indicate LOF depletion, while negative ones indicate excess over expected. Error bars indicate 95% CIs.

379
 380 The pattern of largely neutral dN/dS ratios in control DNMs contrasts with the strong evidence of
 381 negative selection in human single nucleotide polymorphisms (SNPs), particularly in common SNPs.
 382 To shed light on this, we calculated dN/dS ratios for population variants at different allele frequencies
 383 (AF) in the population using data from gnomAD⁵⁹ (**Fig. 4a**). This revealed a decay in dN/dS ratios for
 384 population variants with higher AFs, a pattern consistent with purifying selection operating over many

385 generations on germline SNPs. Altogether, these analyses suggest that the dominant selection force on
386 germline mutations is positive selection during spermatogenesis but negative selection between
387 generations. This dynamic mirrors the well-established contrast between selective forces on cancer
388 mutations and those on germline mutations in populations²⁹.

389

390 We then compared the per-gene enrichment of loss-of-function mutations in sperm to those of
391 population germline variants using the gnomAD loss-of-function z-score. This z-score is a measure of
392 how significantly the observed counts differ from expectations of a germline mutation model. The vast
393 majority of genes in gnomAD have a positive z-score, indicative of depletion of LOFs and the negative
394 selection expected in populations. Of the 31 significant germline selection genes with LOF enrichment
395 in sperm (range 3-fold to 50-fold), 27 were depleted for LOFs in gnomAD. Each of these genes has a
396 disease phenotype, consistent with a model by which these genes are selected for in spermatogenesis
397 but purged from the population by strong negative selection. Interestingly, 4 significant germline
398 selection genes had more loss-of-function mutations than expected in gnomAD: *SMAD6*, *MIB1*, *LZTR1*,
399 and *SSX1* (**Fig. 4b**). The latter 3 of these genes are 3 of the 4 strongest outliers of LOF enrichment of
400 all genes in gnomAD v2 and were given a cautionary outlier label for unexplained LOF enrichment.
401 Our results suggest that the explanation behind their apparent enrichment in gnomAD is that germline
402 positive selection introduces them at a higher rate than for other genes and negative selection against
403 these variants is not strong enough to mask it.

404

405

406

407

408 Discussion

409 We sequenced sperm and blood from healthy men spanning a wide age range to quantify mutation rate
410 and positive selection in the male germline. The observed mutation rates and mutational signatures in
411 sperm were consistent with those from family trio and testis sequencing studies^{7,27,60–63}. We find that,
412 despite sharing the same mutational signatures, SBS1 and SBS5, mutations accumulate at
413 approximately an 8-fold lower rate in sperm compared to blood. This supports our previous observations
414 comparing germline mutations in testes to a wide range of somatic tissues⁷ and emphasises the protected
415 nature of the germline relative to the soma.

416
417 Analysing over 35,000 coding variants from sperm exome-wide, we build on the 13 previously
418 identified germline selection genes, identifying an additional 31 genes under positive selection. This
419 provides a foundational catalogue of genes under selection in the male germline and expands the
420 diversity of pathways and mutational mechanisms linked to selection in this tissue.

421
422 Our findings have significant implications for studies relying on germline mutation models, as they do
423 not currently account for germline positive selection. As shown here, this can lead to inaccurate
424 constraint metrics in population cohorts. Importantly, this bias can also affect the identification of novel
425 disease-causing genes from DNM enrichment tests. For instance, the recent identification of an excess
426 of *de novo* loss-of-function mutations in *MIB1*³⁹ likely reflects germline selection rather than disease
427 association. Loss-of-function variants in *MIB1* are more common in population cohorts than expected
428 for a gene associated with developmental disorders, and carrying one of these variants does not correlate
429 with developmental disorder phenotypes⁶⁴. In principle, adopting case-control tests for DNM
430 enrichment should help to exclude genes under germline selection that are not linked to the disease.
431 However, sufficiently large control trio datasets, as well as close age matching of controls to account
432 for the age dependency of germline selection, will be needed to ensure sufficient statistical power. Until
433 such resources become available, analyses of DNM enrichment in disease should take into account the
434 evidence for germline selection influencing individual genes presented here (**Supplementary Table 3**).

435
436 Unlike the example of *MIB1*, most genes under positive selection during spermatogenesis are known to
437 be associated with severe monogenic disorders with mutation mechanisms under positive selection
438 matching those associated with disease. We demonstrate that this positive selection leads to a striking
439 2-3 fold enrichment in the fraction of sperm carrying a likely disease mutation across the age range
440 studied. As a result, we estimate that 3-4% of sperm from men over 50 carry a likely disease-causing
441 mutation. Somewhat reassuringly, we note that typical paternal ages at conception are younger than the
442 average age of sperm donors studied here, and that the impact of germline positive selection will be

443 correspondingly weaker. Future investigations, including sequencing of sperm from cohorts focused on
444 men under the age of 30 will aid in developing estimates of germline selection strength in this age range.

445

446 A key consideration when interpreting the fraction of disease mutations in sperm is that this fraction
447 may not directly correspond to the rate at which these variants are observed in live births. There are a
448 number of reasons why, for some genes, the transmission rate to live births could be lower than those
449 observed in sperm, including impaired ability of sperm to fertilise an egg, embryonic lethality, or
450 increased pregnancy loss (**Supplementary Note 4**). Future studies, such as sequencing of trio DNMs
451 from large birth cohorts, will be needed to quantify the relationship between the rate of positively
452 selected disease mutations in sperm and disease incidence in populations.

453

454 While up to 3-5% of sperm in middle-aged to elderly individuals harbouring a known driver mutation
455 has a large impact on offspring disease risk, it is on the low end of the spectrum of estimates in
456 proliferative somatic tissues. For instance, in comparable age groups, more than 40% of cells in
457 endometrial and esophageal epithelium carry a driver mutation^{2,5}, whereas only a few percent of cells
458 are estimated to carry a driver mutation in healthy colon or liver^{3,4}. For blood, on average only a few
459 percent of cells carry a driver mutation in middle aged individuals, but some individuals can present
460 large clonal expansions due lack of severe spatial constraints for clonal expansion⁹. While the germline
461 mutation rate is under evolutionary pressure to remain low to prevent detrimental mutations across the
462 genome, it is perhaps under particularly strong pressure to keep driver mutation rate low, as a single
463 germline driver mutation can cause disease in offspring. It is likely that the tubular organisation of the
464 testis provides strong spatial constraints to prevent large clonal expansions, limiting the accumulation
465 of driver mutations in sperm despite the large number of cell divisions required to sustain sperm
466 production. The low relative rate of driver mutations in the male germline is also consistent with the
467 unique aspects of spermatocytic tumours, the tumour type thought to derive from spermatogonial stem
468 cells in the testis. Spermatocytic tumours are rare relative to most somatic tumours and are primarily
469 driven by chromosome aneuploidies⁶⁵ rather than the classical sequential acquisition of driver mutations
470 leading to cancer transformation^{66,67} observed in driver rich somatic tissues.

471

472 The findings of this study provide important insights into the historically underexplored reproductive
473 ageing risks associated with the male germline. This contrasts with the well-established relationship
474 between maternal ageing and reproductive risks, where decreasing oocyte quality with age leads to
475 increased rates of pregnancy loss and aneuploidy⁶⁸. Our results demonstrate that driver mutation
476 accumulation from the male germline's continuous cell proliferation is a substantial risk, though one
477 spread across many genes. As trends toward delayed reproduction continue⁶⁹, it is essential to recognise
478 that both paternal and maternal ageing contribute to elevated risks for offspring, albeit primarily through
479 different biological mechanisms. Future research will refine our understanding of selective pressures

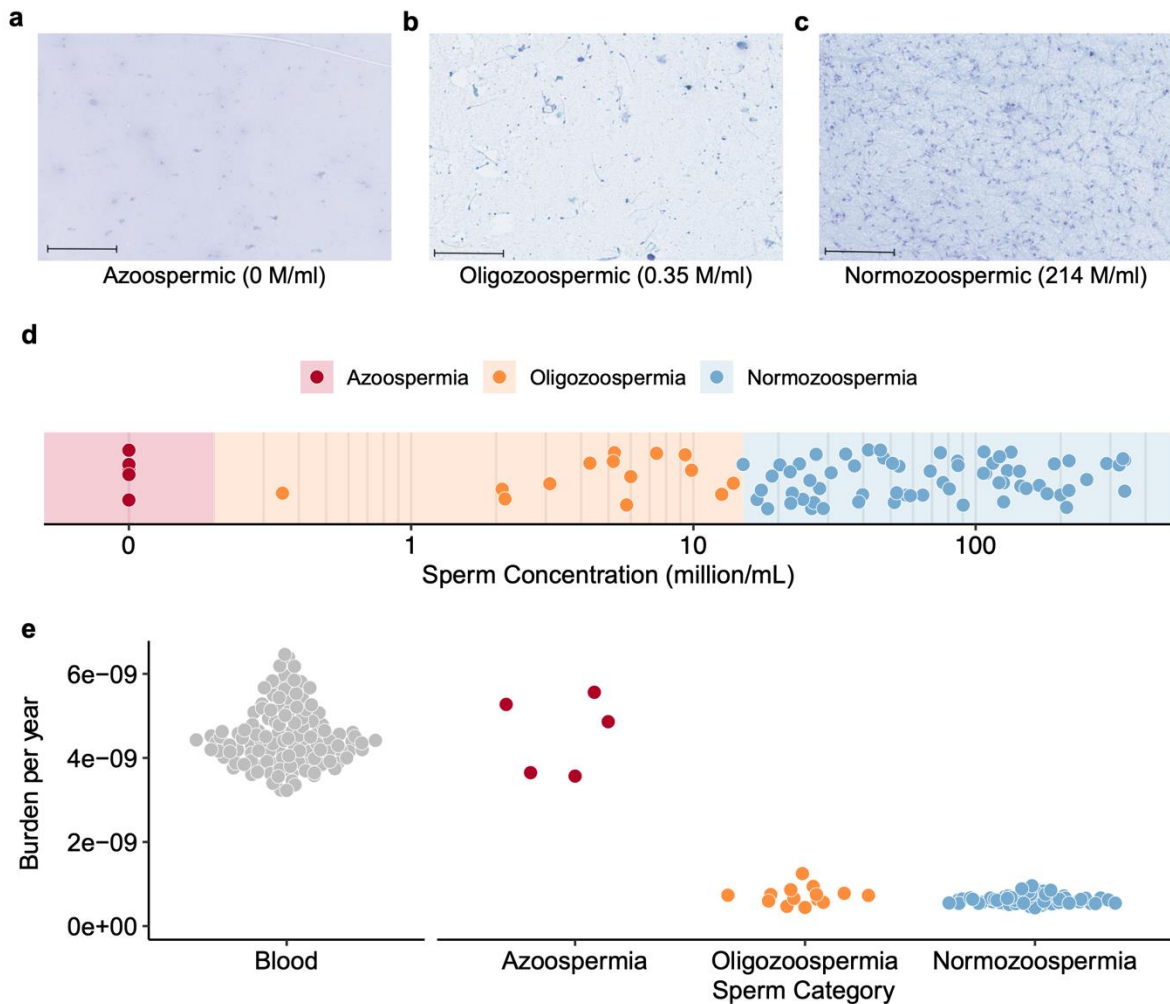
480 and disease risk associated with germline mutation, enhancing our understanding of their implications
481 for human reproduction and health.

482

483

484

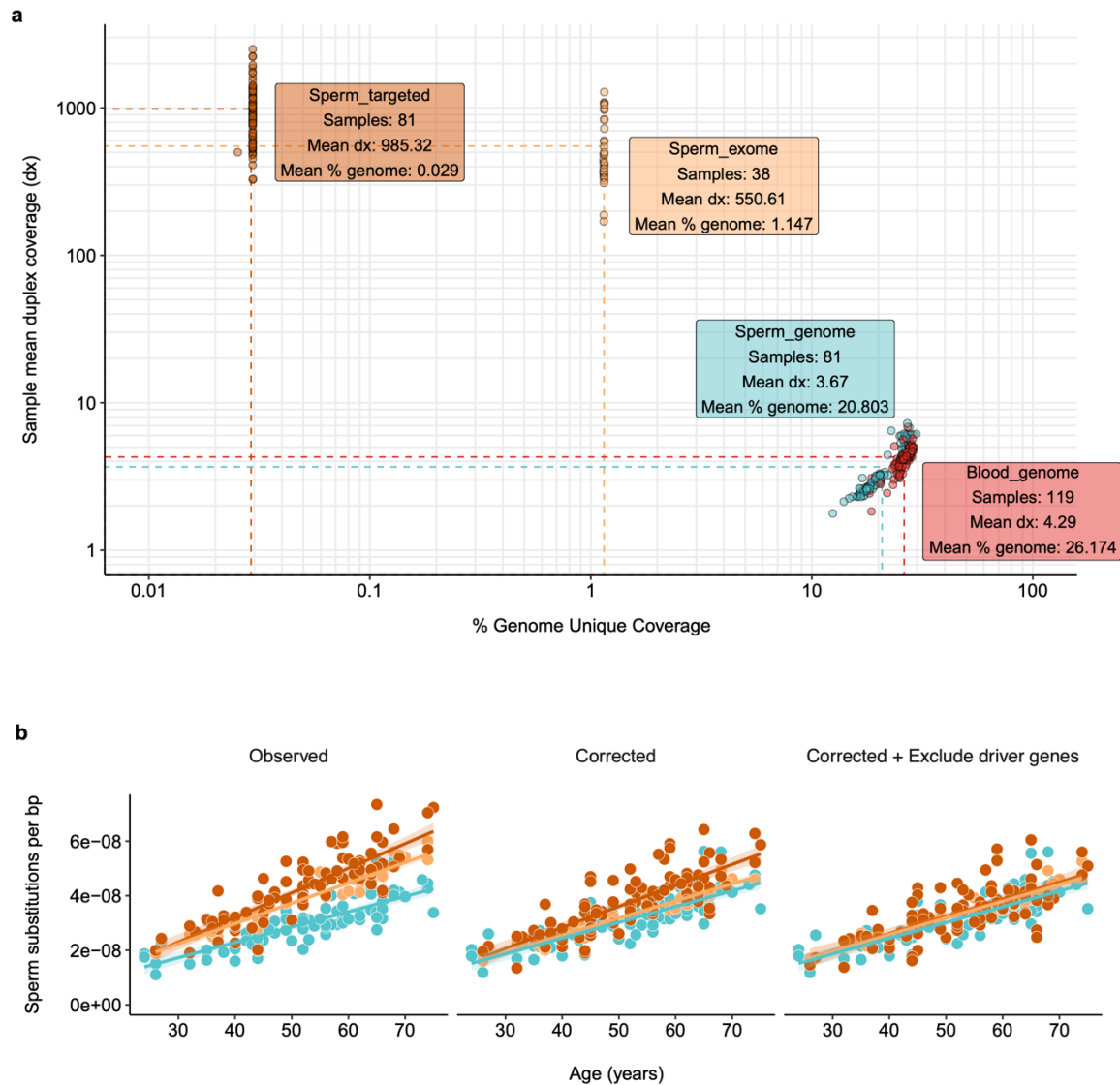
485 **Extended Figures**



486

487 **Extended Data Fig. 1 | Sperm counting**

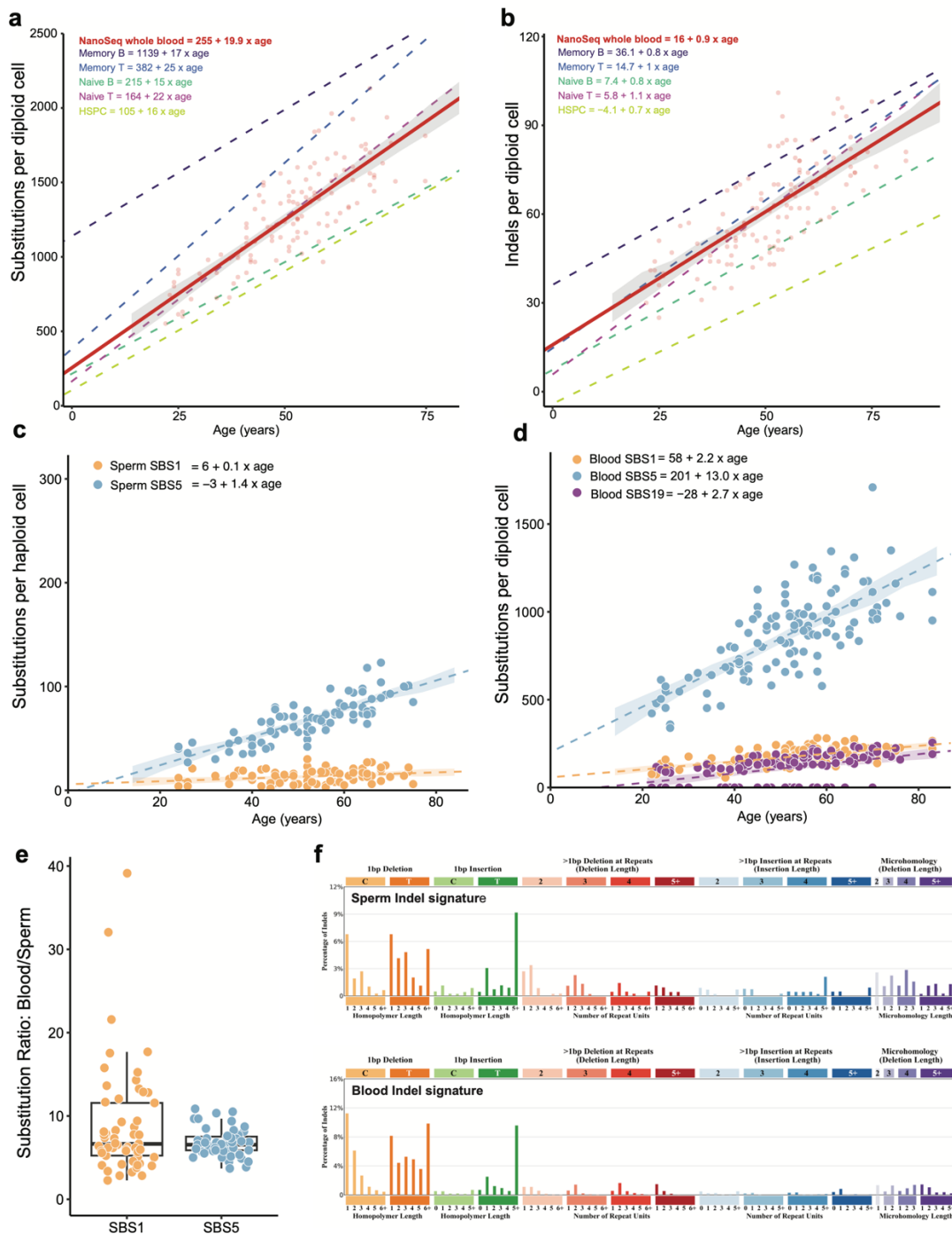
488 **a,b,c**, Slides of Papanicolaou stained semen samples for (a) an azoospermic sample where no sperm cells are
489 visible, (b) an oligozoospermic sample where a small number of sperm samples are visible and (c) a
490 normozoospermic sample where many sperm cells are visible. Sperm concentrations are given for each sample
491 in millions of sperm per ml (M/ml). The black band in the bottom left of each slide photo corresponds to
492 100 μ m. **d**, The distribution of sperm counts on a log scale among semen samples analysed with colour bands
493 indicating the concentration bin of the sample. All samples below 1 million/mL were subsequently excluded. **e**,
494 The distributions of mutation burden per year from blood samples and three categories of sperm samples broken
495 down by sperm concentration.
496



497

498 **Extended Data Fig. 2 | Coverage summary**

499 **a**, Mean duplex coverage (log scale) and percentage of genome covered (log scale) per sample. Panels
500 summarise the mean duplex coverage (dx) and mean percentage of genome covered per NanoSeq type and
501 tissue. **b**, Mutation burden of targeted (dark orange), exome (yellow), and genome (blue) sperm sequenced
502 samples that are observed without correction (left), corrected for trinucleotide composition of covered base pairs
503 relative to the whole genome (middle) or corrected and masked for mutations and coverage in the 44 genes
504 linked to germline positive selection (right). Model fits are linear regressions with 95% CI bands.



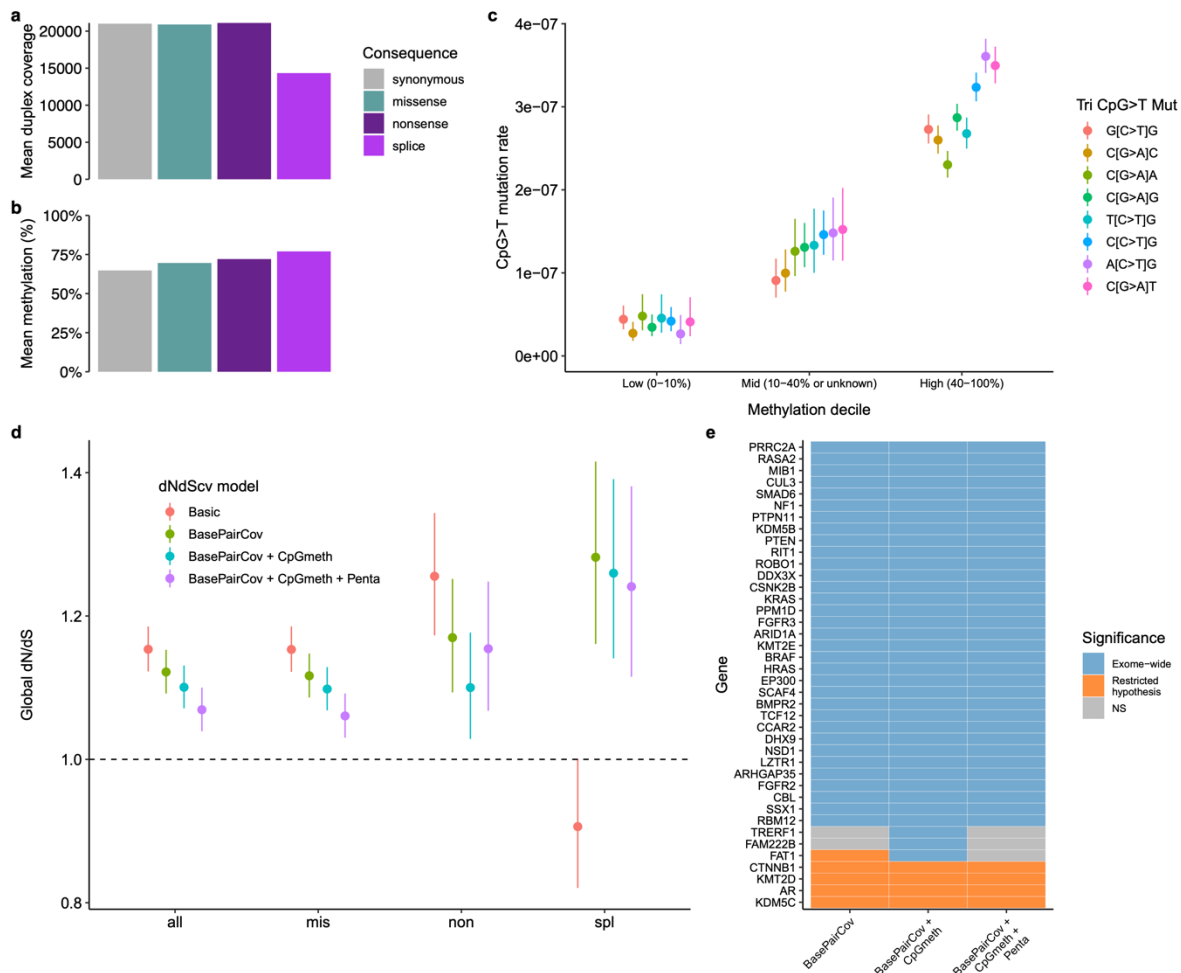
505

506 **Extended Data Fig. 3 | Mutation rates relative to blood cell types and split by signatures**

507 **a,b**, Substitutions (**a**) and indels (**b**) per diploid cell from blood NanoSeq relative to specific blood cell types⁸.

508 **c,d**, Substitutions per haploid cell for sperm (**c**) and diploid cell for blood (**d**) split by signature contributions of
 509 SBS1, SBS5, and SBS19. **a,b,c,d**, Models are linear mixed regressions with 95% CIs calculated by parametric
 510 bootstrapping.

511 **e**, Ratio of age-corrected blood to sperm substitutions per diploid cell per year for mutations
 512 assigned to SBS1 and SBS5. Each dot corresponds to an individual with both a blood and sperm sample and
 513 where individuals had multiple timepoints the mean value of all timepoints in that tissue was used. Box plots
 514 show the interquartile range, median, and 95% confidence interval for the median. **f**, Distribution of indel types
 515 observed in sperm and blood.

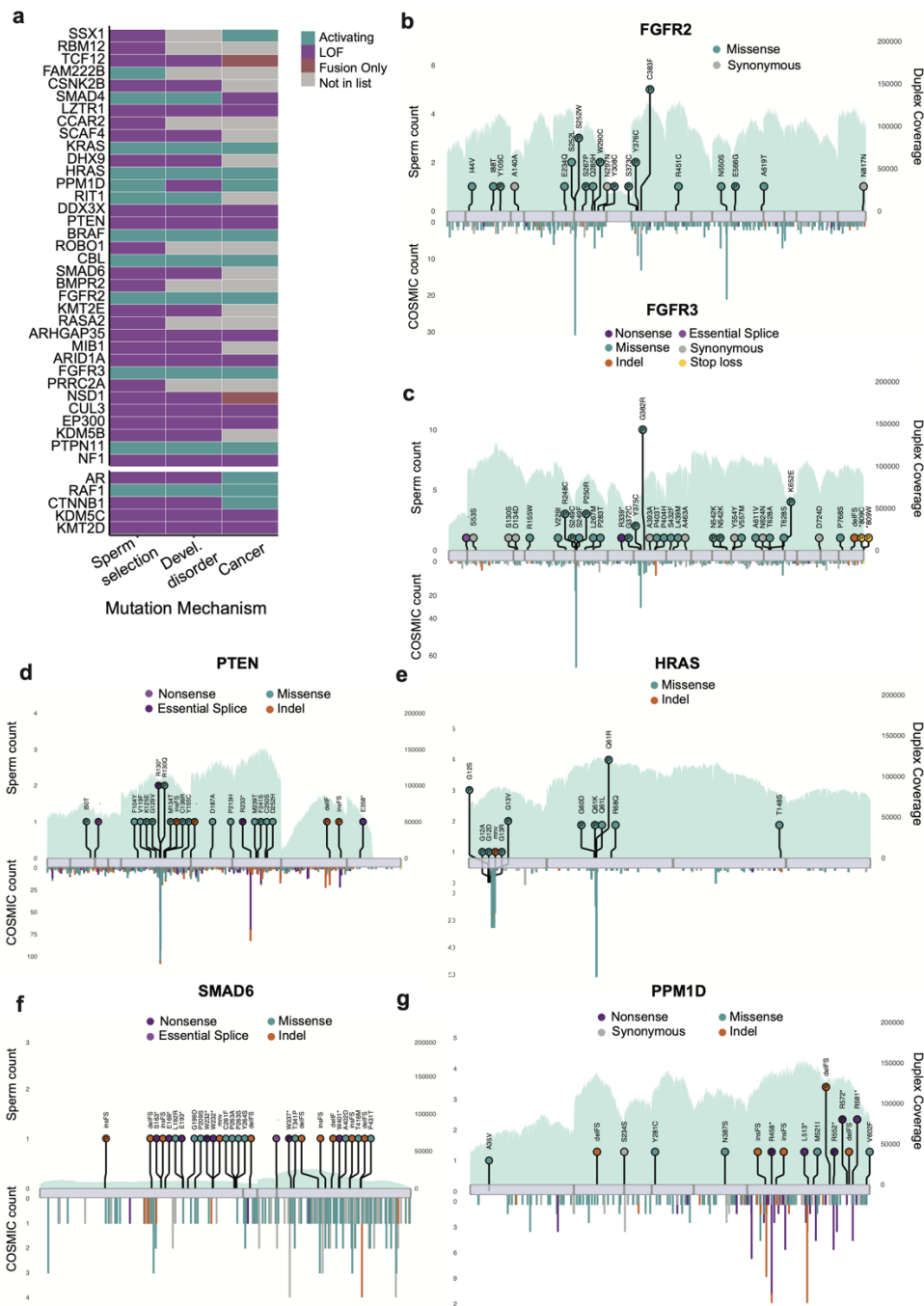


515

516 **Extended Data Fig. 4 | Model selection dN/dS**

517 **a,b**, Mean duplex coverage (**a**) and methylation percentage (**b**) of all base pairs with exome sequencing
 518 coverage split by mutation consequence. **c**, C>T mutation rate at CpG sites in exome sequenced samples split by
 519 methylation bin based on percentage methylated from testis bisulfite sequencing⁷⁰. **d**, Comparison of global
 520 dN/dS values from exome sequenced samples using different modifications to the *dNdScv* algorithm. Categories
 521 are all nonsynonymous mutations, missense, nonsense or essential splice. The basic model excludes genes
 522 which have no coverage but otherwise uses default parameters. Additional models show the impact of adding
 523 corrections for duplex coverage per base pair (BasePairCov), CpG methylation level (CpGmeth), and
 524 pentanucleotide context (Penta). **e**, Comparison of per-gene significance in exome-wide (blue) or restricted
 525 hypothesis (orange) dN/dS tests using the different models. Genes that did not reach significance in either test
 526 are shown in grey. Error bars indicate 95% CIs.

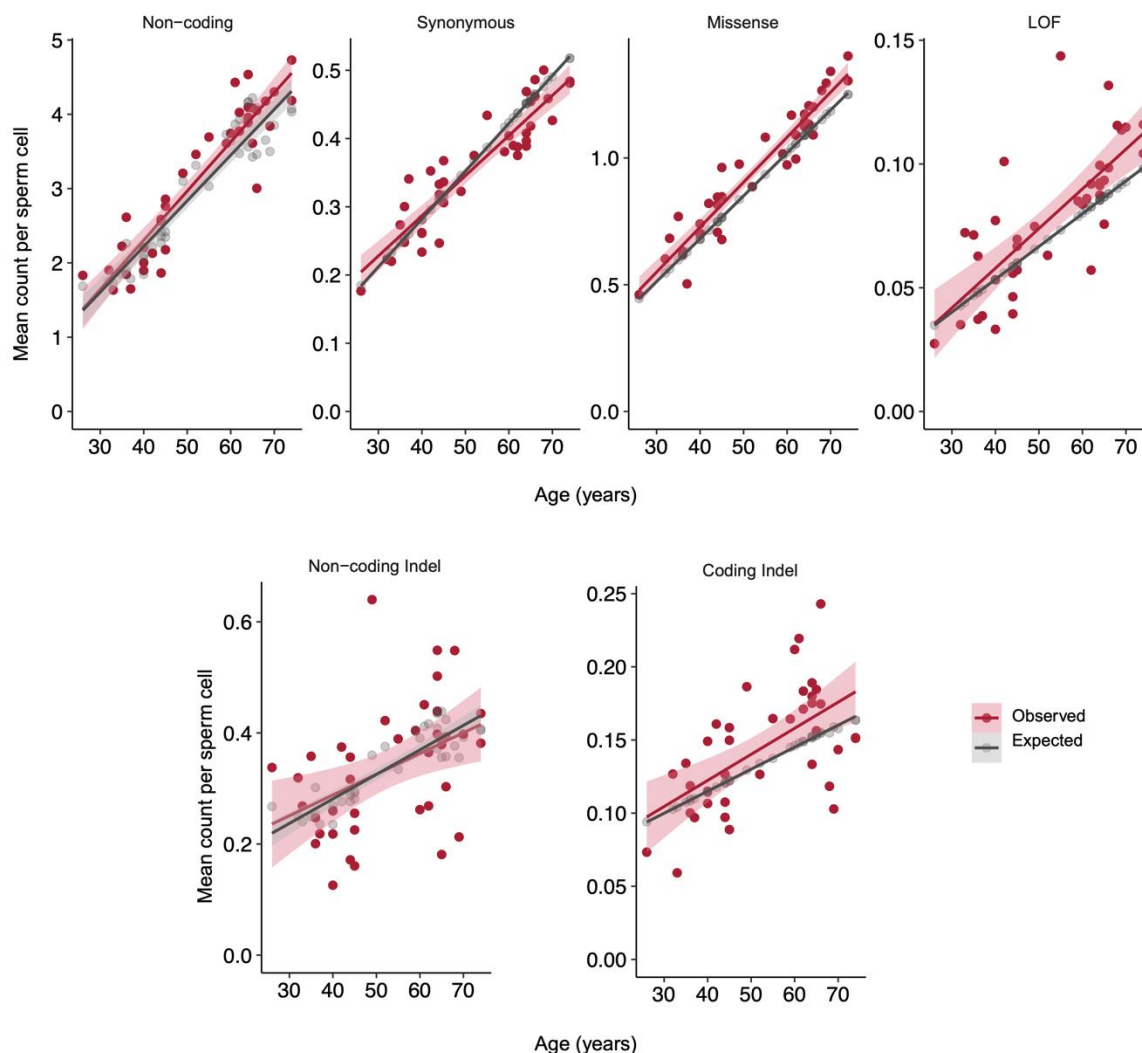
527



528

529 **Extended Data Fig. 5 | Gene mutation mechanisms**

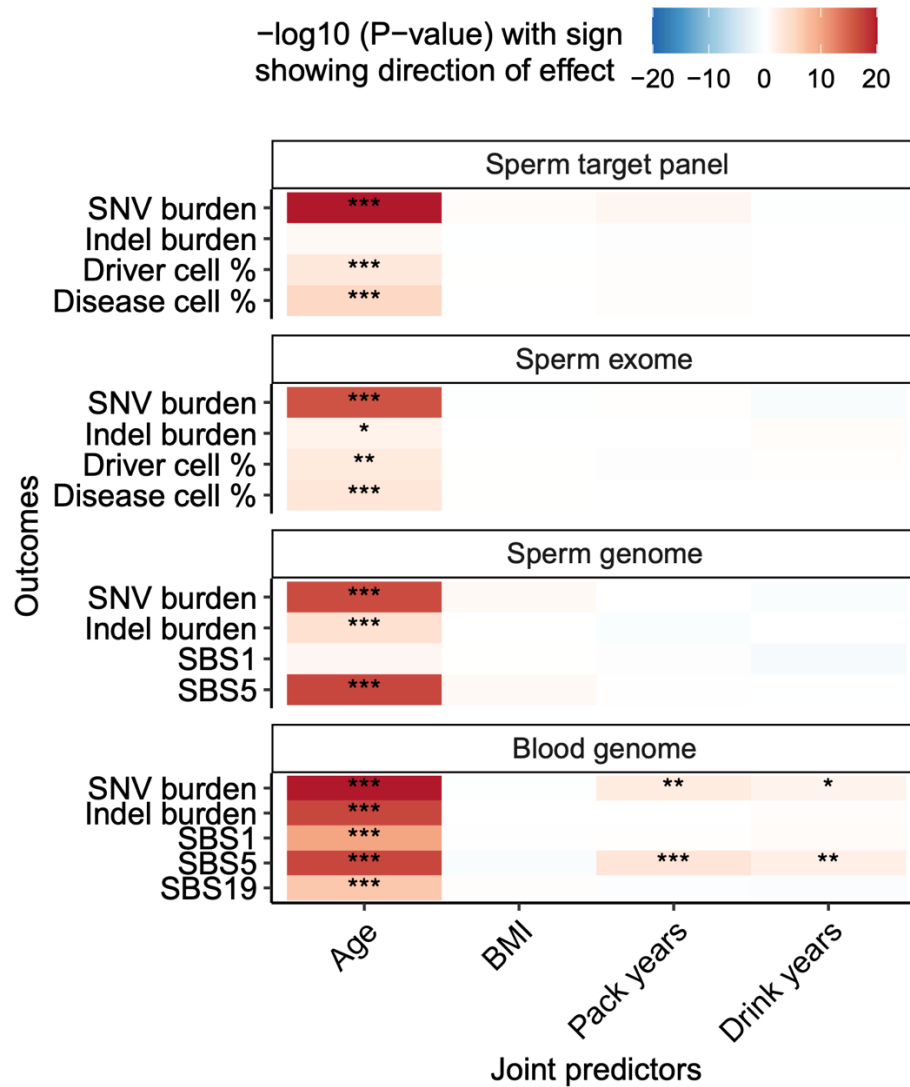
530 **a**, The mutation mechanism assigned to each gene based on the mutation pattern in sperm, developmental
 531 disorders and cancer (**Methods**). **b,c,d,e,f,g**, Observed sperm mutations across the cohort for six illustrative
 532 genes where the height of the “lollipop” represents the number of unique samples with a mutation at that
 533 location and the colour represents its mutation type. Mutations are labelled with their amino acid consequence
 534 for point substitutions or their insertion (ins)/deletion (del) consequence of in frame (IF) or frameshift (FS). A
 535 “P” indicates that the variant is classified as pathogenic/likely pathogenic in ClinVar³⁴. Exons are shown as
 536 purple rectangles and the blue background represents the total duplex coverage across the cohort. Lines below
 537 the gene indicate COSMIC somatic mutations in cancer within that gene³².



538

539 **Extended Data Fig. 6 | Mean variant class count per individual by age**

540 The relationship between age and the mean count of SNVs (non-coding, synonymous, missense, and loss-of-
541 function (nonsense or essential splice)) and indels (non-coding indel and coding indel) per sperm cell. The red
542 points represent the observed values for each individual. The grey line represents the expected mutation count
543 per sperm based on the germline mutation rate model. Error bands indicate 95% CIs of linear regressions.



544

545 **Extended Data Fig. 7 | Phenotype correlations**

546 Correlation of cohort phenotypes to mutation outcome variables, with different sequencing datasets

547 split by facets. Joint predictor glm models used the gaussian family with FDR corrected P values.

548 Asterisks indicate significance level of corrected P value: (* P value >0.01 to <0.05 , ** P value >0.001

549 to <0.01 , *** P value <0.001).

550

chr	pos	ref	mut	gene	aachange	impact	obs	exp	dnds_ratio	qval_exome	qval_RHT
4	1806119	G	A	FGFR3	G382R	Missense	10	0.007750	1290	2.69e-06	2.33e-11
10	123274774	A	G	FGFR2	C383R	Missense	5	0.000803	6226	3.33e-05	2.88e-10
4	1807889	A	G	FGFR3	K652E	Missense	4	0.000225	17816	3.33e-05	2.88e-10
12	112926909	A	G	PTPN11	Q510R	Missense	5	0.000925	5406	4.46e-05	NA
12	112888168	T	G	PTPN11	Y62D	Missense	4	0.000303	13198	6.62e-05	NA
18	48604664	C	T	SMAD4	R496C	Missense	9	0.011300	793	1.81e-04	2.20e-09
12	112888172	A	G	PTPN11	Y63C	Missense	6	0.003010	1994	1.81e-04	2.20e-09
18	48604676	A	G	SMAD4	I500V	Missense	5	0.002230	2244	1.73e-03	2.00e-08
17	27086078	C	T	FAM222B	R300H	Missense	5	0.002320	2157	1.87e-03	NA
22	41572254	T	G	EP300	F1595V	Missense	3	0.000176	17044	2.90e-03	3.58e-08
12	112926890	A	G	PTPN11	M504V	Missense	4	0.000963	4154	2.90e-03	3.58e-08
12	112926887	G	A	PTPN11	G503R	Missense	4	0.000975	4101	2.90e-03	NA
1	155874263	T	C	RIT1	M107V	Missense	5	0.003180	1572	6.09e-03	NA
4	1803571	C	G	FGFR3	P250R	Missense	3	0.000337	8912	1.64e-02	NA
1	155874285	A	C	RIT1	F99L	Missense	3	0.000349	8600	1.70e-02	NA
10	123279677	G	C	FGFR2	S252W	Missense	3	0.000491	6106	4.45e-02	6.85e-07
18	19426998	C	T	MIB1	R769*	Nonsense	4	0.002190	1825	5.05e-02	NA
12	112926270	C	T	PTPN11	T468M	Missense	4	0.003070	1301	1.80e-01	2.81e-06
12	112888199	C	T	PTPN11	A72V	Missense	3	0.000861	3484	2.00e-01	3.00e-06
12	112926851	C	T	PTPN11	P491S	Missense	3	0.000936	3205	2.44e-01	3.50e-06
11	533874	T	C	HRAS	Q61R	Missense	3	0.000959	3127	2.50e-01	3.50e-06
3	12645699	G	A	RAF1	S257L	Missense	4	0.012700	314	1.00e+00	4.61e-04
4	1803564	C	T	FGFR3	R248C	Missense	3	0.005820	516	1.00e+00	6.16e-04
11	534289	C	T	HRAS	G12S	Missense	3	0.010100	297	1.00e+00	2.79e-03
17	29654736	C	T	NF1	R1830C	Missense	3	0.013200	228	1.00e+00	5.53e-03

551

552 **Extended Data Table 1 | Significant SNV hotspots from dN/dS exome-wide and**
553 **restricted hypothesis tests (RHT)**

554

555

556

557 **Methods**

558 **Ethics**

559 This study was carried out under TwinsUK BioBank ethics, approved by North West – Liverpool
560 Central Research Ethics Committee (REC reference 19/NW/0187), IRAS ID 258513 and earlier
561 approvals granted to TwinsUK by the St Thomas’ Hospital Research Ethics Committee, later London
562 – Westminster Research Ethics Committee (REC reference EC04/015).

563

564 **Sample collection**

565 Bulk semen samples were collected or obtained from archival samples with informed consent from 75
566 research participants within the TwinsUK cohort²⁶. Archival whole-blood samples were also obtained
567 from 67 of those men from the TwinsUK BioBank. A total of 104 semen samples spanned an age range
568 of 24-75 years and included 29 men with 2 timepoints separated by a mean of 12.1 years (range 12-13
569 years) and the remaining 46 men with a single timepoint. A total of 133 blood samples were collected
570 at an age range of 22-83 years. There were 11 men with a single blood timepoint, 47 with two
571 timepoints, 8 with three timepoints and 1 with four timepoints. The mean interval between blood
572 timepoints was 8.1 years (range 1-13 years). Within the cohort there were a total of 9 monozygotic
573 (MZ) twins and 3 dizygotic (DZ) twin pairs. Counts of samples, timepoints, and twin pairs which were
574 successfully sequenced and passed analysis quality control thresholds are summarised in
575 **Supplementary Table 1.**

576

577 **Metadata**

578 Metadata for self-reported age, height, weight, ethnicity, twin zygosity, smoking and alcohol
579 consumption were obtained from questionnaires provided by TwinsUK taken periodically. All
580 individuals that provided ethnicity information indicated “white”. BMI was calculated as
581 weight/height². A smoking pack year was defined as 365 packs of cigarettes and total pack years was
582 calculated using the highest estimate across all questionnaires from cigarettes per day or week and total
583 years smoked. Alcohol drink years was calculated from using average weekly alcohol consumption
584 extrapolated to the duration of adult life before sampling (age - 18).

585

586 **DNA extraction**

587 DNA was extracted from sperm samples using the Qiagen QIAamp DNA Blood Mini Kit. Isolation of
588 genomic DNA from sperm; protocol 1 (QA03 Jul-10) was followed with the exceptions of substituting
589 DTT in place of β -mercaptoethanol for Buffer 2 and substituting Buffer EB in place of Buffer AE for
590 elution of DNA.

591 DNA was extracted from whole blood using the Gentra Puregene Blood Kit, following the protocol for
592 10ml of compromised whole blood from the Gentra Puregene Handbook version 06/2011.

593

594 **Targeted gene panels**

595 Three separate Twist Bioscience gene panels were used for targeted NanoSeq sequencing in this study:
596 1) a custom pilot panel of 210 genes; 2) a similar but extended custom panel of 263 genes
597 (**Supplementary Table 2**); 3) a default exome-wide panel of 18,800 genes. The two custom gene panels
598 are highly similar with the extended panel being almost exclusively regions added to the pilot panel.
599 From the 84 samples that underwent targeted sequencing, 13 were sequenced using a pilot panel of 210
600 canonical cancer/somatic driver genes, and all 84 were sequenced using the extended panel of 263
601 genes. Sequencing coverage and mutations were merged from samples sequenced on both targeted
602 panels. The custom panels were designed by gathering sets of published lists of genes implicated as
603 drivers in cancers⁷¹⁻⁷⁴ and somatic tissues^{1,75} as described in (Lawson A.R., Abascal F., P.A. Nicola et
604 al., manuscript submitted for publication).

605

606 **Sequencing and preprocessing of NanoSeq libraries**

607 Restriction-enzyme whole genome NanoSeq libraries were prepared as described in Abascal et al.²³ and
608 subjected to whole genome sequencing at target 20-30x coverage on NovaSeq (Illumina) platforms to
609 generate 150-bp paired-end reads with 9-10 samples per lane. Standard whole genome sequencing of
610 blood (31.7x median coverage) was used to generate matched-normal libraries for both restriction-
611 enzyme NanoSeq blood and sperm.

612

613 Targeted and exome NanoSeq libraries were prepared via sonication and 1-2 rounds of pull down of
614 target sequences as described in (Lawson A.R., Abascal F., P.A. Nicola et al., manuscript submitted for
615 publication). They were then sequenced with NovaSeq (Illumina) platforms to generate 150-bp paired-
616 end reads with 7-8 samples per lane for the targeted panel and 2 lanes per sample for the exome panel.

617

618 **Base calling and filtering**

619 All samples were processed using a Nextflow implementation of the NanoSeq calling pipeline
620 (<https://github.com/cancerit/NanoSeq>). BWA-MEM⁷⁶ was used to align all sequences to the human
621 reference genome (NCBI build37). Restriction-enzyme NanoSeq samples were called with their
622 matched WGS normal and default parameters except for var_b (minimum matched normal reads per
623 strand) of 5 as needed for WGS normals, cov_Q (minimum mapQ to include a duplex read) of 15 and
624 var_n (maximum number of mismatches) of 2.

625

626 For targeted and exome samples we leveraged the high sequencing depth and high polyclonality to
627 exclude variants with VAF > 10% instead of using a matched normal. Default parameters of the calling
628 pipeline except for cov_Q of 30, var_n of 2, var_z (minimum normal coverage) of 25, var_a (minimum
629 AS-XS) of 10, var_v (maximum normal variant allele frequency (VAF)) of 0.1, and indel_v (maximum
630 normal VAF) of 0.1. Post variant calling, we further filtered variants to those below 1% VAF and those
631 below 10% duplex VAF as variants above these cutoffs were highly enriched for mapping artefacts,
632 particularly for indels. No excluded variants from these additional VAF thresholds were found to be
633 likely driver or ClinVar pathogenic variants; all exclusions were inspected to confirm this.

634

635 A set of common germline variants from dbsnp⁷⁷ and a custom set of known artifactual call sites in
636 NanoSeq datasets were masked for coverage and variant calls as previously described²³.

637

638 **Assessing DNA contamination**

639 The single-molecule accuracy of the duplex sequencing method NanoSeq allows sequencing of
640 polyclonal cell types such as sperm, but also renders mutation calls sensitive to a) non-target cell-type
641 contamination and b) contamination of foreign DNA. Non-target cell-type contamination was evaluated
642 using manual cell counting of semen samples, resulting in the exclusion of 6 samples with sperm count
643 < 1 million sperm/mL. Sperm counting methods and analysis are detailed in **Supplementary Note 1**.

644

645 Foreign DNA contamination in whole genome NanoSeq samples was assessed using verifyBamID⁷⁸,
646 which checks whether reads in a BAM file match previous genotypes for a specific sample, with higher
647 values indicating more contamination. Three blood whole genome NanoSeq samples were excluded
648 based on a verifyBamID alpha value above the suggested cutoff of 0.005²³. In sperm, we found that
649 several samples had outlier mutation burdens with verifyBamID values just below the 0.005 cutoff.
650 This is logical, as sperm has a much lower mutation rate compared to somatic tissues, for which the
651 recommended cutoff was designed. Consequently, sperm samples are more sensitive to low levels of
652 contamination. To account for this, we adjusted the verifyBamID alpha threshold for sperm to a more
653 stringent level of 0.002, resulting in the exclusion of 3 samples on this criterion.

654

655 When assessing foreign DNA contamination in targeted and exome samples we found that 9 targeted
656 and 3 exome samples had verifyBamID values above > 0.002, slight outlier mutation burdens, and high
657 ratio of SNP masked variants to passed variants (4-fold to 16-fold more masked variants). Upon further
658 investigation we found that all samples exceeding verifyBamID thresholds were processed in the same
659 sequencing batch and that this contamination could be explained by inherited germline variants of other
660 samples within that same batch. This suggests that a small amount of cross-contamination may have
661 occurred during sample preparation or sequencing steps. In order to remove contaminant germline

662 mutations we performed an *in silico* decontamination as previously described²³. This involved calling
663 germline variants from all targeted and exome samples using bcftools mpileup⁷⁹ at sites where there
664 were >10 reads and a mutation call with VAF > 0.3. All such sites were subsequently masked across all
665 samples for both mutation calls and coverage, essentially extending the default common SNP mask to
666 also include rare inherited variants across the cohort. This resulted in all samples previously identified
667 as contaminated having mutation burdens consistent with their age and all having a ratio of masked to
668 passed variants <0.1, and were thus retained for analysis.

669

670 **Corrected mutation burdens**

671 Given that mutation rates are strongly influenced by trinucleotide composition, it is important to
672 consider differences in sequence composition when comparing mutation rates in datasets that target
673 different regions of the genome. For instance, it is known that coding regions such as those in NanoSeq
674 target panels, are biased towards a higher mutation rate partially due to a higher GC density than non-
675 coding regions⁸⁰ which make up the majority of sequenced regions in whole genome NanoSeq datasets.
676 To correct for this effect, in each sperm NanoSeq dataset we generated a corrected mutation burden
677 relative to the full genome trinucleotide frequencies as described previously²³.

678

679 **Comparison of NanoSeq and WGS burden estimates**

680 In order to compare whole genome NanoSeq mutation burdens to mutation burden from standard whole
681 genome sequencing (WGS) we multiplied the corrected mutation burden estimates described in the
682 previous section by the genome size per cell type. We assumed 2,861,326,455 mappable base pairs in
683 a haploid cell for germline datasets and the diploid equivalent of 5,722,652,910 base pairs for blood.

684

685 External datasets for comparison to NanoSeq results were processed in order to achieve comparable
686 burden estimates. For testis WGS samples⁷⁰ we implemented the method described in Abascal et al.²³
687 that restricts analysis to regions with high coverage (20+ reads) that overlap with NanoSeq covered
688 regions. Additionally, we corrected for differences in trinucleotide background frequencies relative to
689 the full genome as described in the previous section. For trio paternally phased DNMs from standard
690 sequencing, as a callable genome size per sample following thorough filtering was available, we
691 generated the mutation per cell estimate by multiplying the paternally phased DNM count by the ratio
692 of the sample's callable genome to total genome size. For comparison to cell types in blood we
693 compared directly to the published mutation burden regressions⁸.

694

695 **Mutation burden regressions**

696 To investigate the relationship between age and mutation burdens, we performed linear mixed-effects
697 regression analyses. For each tissue and mutation type where a regression was performed, the model

698 was constructed using the `lmer` function from the `lme4` package⁸¹ in R. Each model included age at
699 sampling as a fixed effect and a random slope for each individual to account for multiple timepoint
700 samples, specified as:

701

702
$$\text{lmer}(\text{burden} \sim \text{age} + (\text{age} - 1 | \text{indiv_ID}), \text{REML}=\text{F})$$

703

704 The 95% confidence intervals (CIs) for regression lines were calculated through bootstrapping by
705 simulating prediction intervals. For each model, we generated 1000 bootstrap samples. Predictions and
706 their associated standard errors were calculated for a sequence of ages from 14 to 84 years. The 95%
707 CIs were then derived by determining the range within which 95% of the bootstrap sample predictions
708 fell.

709

710 **Mutational signature analysis**

711 We extracted *de novo* mutational signatures using Hierarchical Dirichlet Process (HDP;
712 <https://github.com/nicolaroberts/hdp>), which is based on the Bayesian Hierarchical Dirichlet process.
713 HDP was run with double hierarchy: (1) individual ID and (2) tissue types (either blood or sperm), and
714 without the Catalogue Of Somatic Mutations In Cancer (COSMIC) reference signatures⁸² (v3.3) as
715 priors, on the mutation matrices. The number of mutations were normalised for the tri-nucleotide
716 context abundance specific for each sample relative to the full genome. Both clustering hyper-
717 parameters, beta and alpha, were set to one. The Gibbs samples ran for 30,000 burn-in iterations
718 (parameter “burnin”), with a spacing of 200 iterations (parameter “space”), from which 100 iterations
719 were collected (parameter “n”). After each Gibbs iteration, three iterations of concentration parameters
720 were conducted (parameter “cpiter”). Two components were extracted as *de novo* mutational signatures
721 which were further reconstructed and decomposed into known COSMICv3.3 SBS signatures using
722 SigProfilerAssignment (<https://github.com/AlexandrovLab/SigProfilerAssignment>). As a result, three
723 COSMIC signatures: SBS1, SBS5, and SBS19, were reported.

724 **Quantifying selection with dN/dS**

725 To examine genes under positive selection and quantify global selection we used the *dNdScv*
726 algorithm²⁹. This algorithm was extended using base pair level duplex coverage, methylation level and
727 pentanucleotide context to capture more complex context dependent mutational biases, and to achieve
728 more accuracy for our selection analysis. Detailed methods for input mutations, model selection and
729 evaluation, site dN/dS tests, driver mutation estimation, and gene set enrichment are described in
730 **Supplementary Note 3.**

731 **Gene disease and mechanism annotation**

732 Positively selected genes were annotated with monoallelic disease consequences using the 2024-02-29
733 release of the Development Disorder Genotype - Phenotype Database (DDG2P)³³ and the 2024-06-21
734 release of the Online Mendelian Inheritance in Man (OMIM) database⁴⁵. OMIM annotations related to
735 somatic disease, complex disease, or tentative disease associations were excluded.

736

737 Genes were also annotated for their potential mutation mechanism observed in sperm and
738 cancer/developmental disorders when available. In sperm, genes were labelled as loss-of-function if
739 they had nominal enrichment of nonsense+splice variants and/or indel variants ($\text{p trunc_cv} < 0.1$ |
740 $\text{p ind_cv} < 0.1$) and 2+ loss-of-function mutations. There were two exceptions to this where genes met
741 these thresholds but were labelled as activating due to having a restricted repertoire of loss of function
742 mutations that are known to be oncogenic in cancers: *CBL* (LOFs in and downstream of the RING zinc
743 finger domain)⁸³ and *PPM1D* (LOFs in final two exons)⁸⁴. All other genes had missense enrichment
744 only and were labelled as activating. The mechanism in cancer was defined by using the 'Role in Cancer'
745 field of the COSMIC cancer gene census v99³² where 'oncogene' was labelled as activating and 'tumour
746 suppressor gene' as loss of function. Annotations of a fusion mechanism were not displayed except for
747 genes which had neither an oncogene, nor a tumour suppressor annotation which were labelled as
748 'fusion only'. The developmental disorder mechanism was defined by using the variation consequence
749 field of DDG2P where 'restricted repertoire of mutations;activating' was labelled as activating and
750 'loss_of_function_variant' was labelled as loss of function.

751 **Gene mutation plots**

752 The "lollipop" gene mutation plots were created with a coordinate system where the 1 was the first
753 coding base of the GRCh37-GencodeV18+Appris⁸⁵ transcript of the gene. The data sources included
754 protein domains from UniProt⁸⁶, somatic mutations from the exome and genome wide screens of the
755 COSMIC³² (v99), ClinVar release 2024.07.01³⁴ pathogenic annotation, per base pair cohort wide
756 (targeted + exome) NanoSeq coverage, sperm mutation count (number of independent individuals with
757 a mutation) and mutation consequence and amino acid change annotated by the dNdScv algorithm²⁹.
758 These data were plotted with code modified from the lollipop function in the trackViewer R package⁸⁷.

759 **Variant annotation**

760 Variants were annotated using Ensembl's Variant Effect Predictor (VEP)⁸⁸ with added custom
761 annotations of mutation context, ClinVar release 2024.07.01³⁴, Combined Annotation Dependent
762 Depletion (CADD) version GRCh37-v1.6⁵³ and average methylation level in testis. Methylation data
763 was obtained from whole genome shotgun bisulfite sequencing methylation data from a 37 year old
764 (ENCFF638QVP) and 54 year old (ENCFF715DMX) male testis from the ENCODE project⁷⁰. The

765 average methylation level was calculated by selecting CpG sites with coverage of 3 or more and
766 averaging the percent of sites methylated between the two samples.

767

768 Variants were annotated as likely monoallelic disease-causing mutations if they met at least one of two
769 criteria: 1) Reported in ClinVar as pathogenic, likely pathogenic, or if they were reported as
770 ‘Conflicting_classifications_of_pathogenicity’ where the conflict was between reports of
771 pathogenic/likely pathogenic and ‘Uncertain_significance’ with no reports of benign or likely benign
772 and not specified as a recessive condition or 2) Were a highly damaging variant in a high confidence
773 monoallelic developmental disorder genes from DDG2P³³. Genes met criteria of a) allelic requirement
774 being monoallelic_autosomal, monoallelic_X_hem, monoallelic_X_het, or mitochondrial, b)
775 confidence in strong, definitive, or moderate and c) a mutation consequence of ‘absent gene product’.
776 Highly damaging was defined as being a ‘HIGH’ impact variant in VEP annotation (frameshift
777 splice_acceptor, splice_donor, start_lost, stop_gained, or stop_lost) or a missense variant with CADD⁵³
778 score >30 (top 0.1% damaging).

779

780 Variants were defined as a likely driver if they met the ‘highly damaging’ criteria defined above in a
781 significant germline selection gene with loss-of-function mutation enrichment or if they were in one of
782 the 24 significant mutation hotspots. This resulted in 320 variants being labelled as likely drivers in
783 exome samples.

784 **Cell fraction mutation estimates**

785 To calculate the mean count of synonymous, missense or loss-of-function, or pathogenic mutations per
786 sperm cell we summed the duplex VAFs of all variants in that class. For example, if an individual had
787 three synonymous mutations each observed once with a duplex coverage of 100 at each of those sites,
788 each of those variants would have a duplex VAF of $1/100 = 0.01$. The sum of VAFs in this example
789 would then be 0.03 and this would then be reported as the estimate for the mean count of synonymous
790 variants per sperm cell for that individual. At low fractions such as 0.03, the mean count per cell is
791 approximately equivalent to the percentage of sperm with this mutation class (3%) and thus the driver
792 and disease mutations are reported as percentage estimates. At higher fractions (e.g. mean count > 1)
793 the fractions are not equivalent to percentage as many cells will have multiple variants of that class and
794 thus the estimates are reported as mean count.

795

796 Expected mean counts for SNVs were generated by annotating each possible substitution at each
797 covered site with an expected number of mutations per sample as given by $\text{expCountSNV} =$
798 $\text{context_mut_rate} * \text{duplex_coverage} * \text{age_correction}$. The context_mut_rate was given by the 208
799 basePairCov + cpGMeth trinucleotide mutation model estimates for that trinucleotide+methylation

800 mutation context (**Supplementary Note 3; Supplementary Table 6**). Duplex coverage is the exact
801 duplex coverage at that site for that sample. The `age_correction` parameter was given to normalise the
802 mutation model estimates (derived from all exome samples) to the mutation rate of that sample based
803 on age. Specifically, we fit a linear model to the mutation burden vs age of the exome samples and used
804 this to generate a predicted mutation rate for each sample based on age. The per sample corrected
805 parameter was calculated as the age predicted mutation burden divided by the mean mutation rate of all
806 exome samples ($3.89e-08$). The resulting corrections spanned from 0.50 (youngest sample) to 1.42
807 (oldest samples). The expected indel mutation rate was calculated in the same way, except with a single
808 mutation rate parameter (indels/bp) $\text{expCountIndel} =$
809 $\text{indel_mut_rate} * \text{duplex_coverage} * \text{age_correction}$. The expected mean count was then calculated for
810 each category (e.g. synonymous, likely disease) by summing the expected values for each SNV and/or
811 indel base pair matching the relevant annotation. As background for possible ClinVar pathogenic
812 variants we only considered indels of size 21 bp or less, the largest detected indel in the dataset.
813 Regressions were fit with either linear models or generalised linear models (*glm* in R) with family =
814 quasibinomial.

815 **Regression analysis**

816 We tested for associations between mutation outcome variables from sperm genome, sperm exome,
817 sperm targeted and blood genome NanoSeq data and the phenotype predictor variables of BMI, smoking
818 pack years, and alcohol drink years. These tests were performed using a gaussian family generalised
819 linear regression in R. For each mutation outcome variable the test took the form of:
820 $\text{glm}(\text{mutationOutcome} \sim \text{age_at_sampling} + \text{BMI} + \text{pack_years} + \text{drinkYears}, \text{family} = \text{"gaussian"})$.

821
822 The mutation outcome variables examined were SNV and indel burden from all 4 sequencing datasets,
823 SBS1 and SBS5 count from sperm genomes, SBS1, SBS5, and SBS19 from blood genomes, and likely
824 disease cell fraction and likely driver cell fraction from sperm targeted and sperm exomes. The
825 significance of each predictor was assessed from the model's summary output coefficients, and p-values
826 were adjusted for 68 total tests using the false discovery rate method.

827

828

829

830

831

832

833

834 **Data availability**

835 Raw sequencing data are available on the European Genome–Phenome Archive under accession
836 number X. Additional individual-level data are not permitted to be publicly shared or deposited due to
837 the original consent given at the time of data collection, where access to these data is subject to
838 governance oversight. All data access requests are overseen by the TwinsUK Resource Executive
839 Committee (TREC). For information on access to these genotype and phenotype data and how to apply,
840 see <https://twinsuk.ac.uk/resources-for-researchers/access-our-data/>.

841

842 **Code availability**

843 All scripts are available on github at <https://github.com/mattnev17/spermPositiveSelectionManuscript>.

844

845 **Acknowledgements**

846 We thank the TwinsUK research volunteers for participating in the study. We are grateful to Laura
847 O’Neill, Calli Latimer and all of the CASM Support team at the Wellcome Sanger Institute for their
848 assistance. We thank Charlotte Seymour, Elisa Ferraro, and Chris White from Cambridge IVF for
849 training and guidance on sperm counting.

850

851 **Funding**

852 This research is supported by core funding from Wellcome Trust. R.R. is funded by Cancer Research
853 UK (C66259/A27114) and Medical Research Council (MR/W025353/1). TwinsUK is funded by the
854 Wellcome Trust, Medical Research Council, Versus Arthritis, European Union Horizon 2020, Chronic
855 Disease Research Foundation (CDRF), Zoe Ltd, the National Institute for Health and Care Research
856 (NIHR) Clinical Research Network (CRN) and Biomedical Research Centre based at Guy’s and St
857 Thomas’ NHS Foundation Trust in partnership with King’s College London.

858

859 **Author Contributions**

860 M.D.C.N., M.E.H. and R.R. wrote the manuscript; all authors reviewed and edited the manuscript.
861 M.E.H., and R.R. supervised the project. M.P.G., S.W., K.S., and R.R. led sample procurement.
862 M.D.C.N., T.B., and P.A.N. conducted sperm counting. A.R.J.L., P.J.C., K.R., S.V.L., S.W., and I.M.
863 contributed to sample sequencing implementation. M.D.C.N., led the analysis of the data with help from
864 A.R.J.L., R.S., F.A., M.H.P., A.C., P.A.N., A.B.O., R.E.A., M.R.S., P.J.C., I.M., M.E.H., and R.R.

865

866 **Competing Interests**

867 I.M., M.R.S., and P.J.C. are co-founders, shareholders, and consultants for Quotient Therapeutics Ltd.

868 R.E.A. is an employee of Quotient Therapeutics Ltd. M.E.H. is a co-founder of, consultant to and holds
869 shares in Congenica, a genetics diagnostic company.

870

871

872

873

874 **References**

- 875 1. Martincorena, I. *et al.* High burden and pervasive positive selection of somatic mutations in normal
876 human skin. *Science* **348**, 880–886 (2015).
- 877 2. Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age. *Science* **362**,
878 911–917 (2018).
- 879 3. Brunner, S. F. *et al.* Somatic mutations and clonal dynamics in healthy and cirrhotic human liver.
880 *Nature* **574**, 538–542 (2019).
- 881 4. Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial cells. *Nature*
882 **574**, 532–537 (2019).
- 883 5. Moore, L. *et al.* The mutational landscape of normal human endometrial epithelium. *Nature* **580**,
884 640–646 (2020).
- 885 6. Lawson, A. R. J. *et al.* Extensive heterogeneity in somatic mutation and selection in the human
886 bladder. *Science* **370**, 75–82 (2020).
- 887 7. Moore, L. *et al.* The mutational landscape of human somatic and germline cells. *Nature* **597**, 381–
888 386 (2021).
- 889 8. Machado, H. E. *et al.* Diverse mutational landscapes in human lymphocytes. *Nature* **608**, 724–732
890 (2022).
- 891 9. Mitchell, E. *et al.* Clonal dynamics of haematopoiesis across the human lifespan. *Nature* **606**, 343–
892 350 (2022).
- 893 10. Bernstein, N. *et al.* Analysis of somatic mutations in whole blood from 200,618 individuals
894 identifies pervasive positive selection and novel drivers of clonal hematopoiesis. *Nat. Genet.* 1–9
895 (2024) doi:10.1038/s41588-024-01755-1.
- 896 11. Heller, C. H. & Clermont, Y. Kinetics of the germinal epithelium in man. *Recent Prog. Horm.*
897 *Res.* **20**, 545–575 (1964).
- 898 12. Neto, F. T. L., Bach, P. V., Najari, B. B., Li, P. S. & Goldstein, M. Spermatogenesis in
899 humans and its affecting factors. *Semin. Cell Dev. Biol.* **59**, 10–26 (2016).
- 900 13. Goriely, A., McVean, G. A. T., Røjmyr, M., Ingemarsson, B. & Wilkie, A. O. M. Evidence
901 for Selective Advantage of Pathogenic FGFR2 Mutations in the Male Germ Line. *Science* **301**,

- 902 643–646 (2003).
- 903 14. Choi, S.-K., Yoon, S.-R., Calabrese, P. & Arnheim, N. Positive Selection for New Disease
904 Mutations in the Human Germline: Evidence from the Heritable Cancer Syndrome Multiple
905 Endocrine Neoplasia Type 2B. *PLOS Genet.* **8**, e1002420 (2012).
- 906 15. Maher, G. J., Goriely, A. & Wilkie, A. O. M. Cellular evidence for selfish spermatogonial
907 selection in aged human testes. *Andrology* **2**, 304–314 (2014).
- 908 16. Maher, G. J. *et al.* Visualizing the origins of selfish de novo mutations in individual
909 seminiferous tubules of human testes. *Proc. Natl. Acad. Sci.* **113**, 2454–2459 (2016).
- 910 17. Eboeime, J. *et al.* Germline selection of PTPN11 (HGNC:9644) variants make a major
911 contribution to both Noonan syndrome’s high birth rate and the transmission of sporadic cancer
912 variants resulting in fetal abnormality. *Hum. Mutat.* **43**, 2205–2221 (2022).
- 913 18. Wood, K. A. & Goriely, A. The impact of paternal age on new mutations and disease in the
914 next generation. *Fertil. Steril.* (2022) doi:10.1016/j.fertnstert.2022.10.017.
- 915 19. Goriely, A. & Wilkie, A. O. M. Paternal age effect mutations and selfish spermatogonial
916 selection: causes and consequences for human disease. *Am. J. Hum. Genet.* **90**, 175–200 (2012).
- 917 20. Kennedy, S. R. *et al.* Detecting ultralow-frequency mutations by Duplex Sequencing. *Nat.*
918 *Protoc.* **9**, 2586–2606 (2014).
- 919 21. Schmitt, M. W. *et al.* Detection of ultra-rare mutations by next-generation sequencing. *Proc.*
920 *Natl. Acad. Sci.* **109**, 14508–14513 (2012).
- 921 22. Hoang, M. L. *et al.* Genome-wide quantification of rare somatic mutations in normal human
922 tissues using massively parallel sequencing. *Proc. Natl. Acad. Sci.* **113**, 9846–9851 (2016).
- 923 23. Abascal, F. *et al.* Somatic mutation landscapes at single-molecule resolution. *Nature* **593**,
924 405–410 (2021).
- 925 24. Liu, M. H. *et al.* DNA mismatch and damage patterns revealed by single-molecule
926 sequencing. *Nature* 1–10 (2024) doi:10.1038/s41586-024-07532-8.
- 927 25. Kunisaki, J. *et al.* Sperm from infertile, oligozoospermic men have elevated mutation rates.
928 Preprint at <https://doi.org/10.1101/2024.08.22.24312232> (2024).
- 929 26. Verdi, S. *et al.* TwinsUK: The UK Adult Twin Registry Update. *Twin Res. Hum. Genet.* **22**,

- 930 523–529 (2019).
- 931 27. Sasani, T. A. *et al.* Large, three-generation human families reveal post-zygotic mosaicism and
932 variability in germline mutation accumulation. *eLife* **8**, e46922 (2019).
- 933 28. Campbell, P. *et al.* Prolonged persistence of mutagenic DNA lesions in stem cells. Preprint at
934 <https://doi.org/10.21203/rs.3.rs-3610927/v1> (2023).
- 935 29. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell*
936 **171**, 1029-1041.e21 (2017).
- 937 30. Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation rate
938 heterogeneity in humans. *Nat. Commun.* **9**, 3753 (2018).
- 939 31. Xia, B. *et al.* Widespread Transcriptional Scanning in the Testis Modulates Gene Evolution
940 Rates. *Cell* **180**, 248-262.e21 (2020).
- 941 32. Sondka, Z. *et al.* The COSMIC Cancer Gene Census: describing genetic dysfunction across
942 all human cancers. *Nat. Rev. Cancer* **18**, 696–705 (2018).
- 943 33. Thormann, A. *et al.* Flexible and scalable diagnostic filtering of genomic variants using G2P
944 with Ensembl VEP. *Nat. Commun.* **10**, 2373 (2019).
- 945 34. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting
946 evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
- 947 35. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new
948 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361
949 (2017).
- 950 36. Digre, A. & Lindskog, C. The Human Protein Atlas—Spatial localization of the human
951 proteome in health and disease. *Protein Sci. Publ. Protein Soc.* **30**, 218–233 (2021).
- 952 37. Imielinski, M., Guo, G. & Meyerson, M. Insertions and Deletions Target Lineage-Defining
953 Genes in Human Cancers. *Cell* **168**, 460-472.e14 (2017).
- 954 38. Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer whole genomes.
955 *Nature* **578**, 102–111 (2020).
- 956 39. Kaplanis, J. *et al.* Evidence for 28 genetic disorders discovered by combining healthcare and
957 research data. *Nature* **586**, 757–762 (2020).

- 958 40. Wood, K. A. *et al.* *SMAD4* mutations causing Myhre syndrome are under positive selection in
959 the male germline. *Am. J. Hum. Genet.* (2024) doi:10.1016/j.ajhg.2024.07.006.
- 960 41. Waite, K. A. & Eng, C. From developmental disorder to heritable cancer: it's all in the
961 BMP/TGF- β family. *Nat. Rev. Genet.* **4**, 763–773 (2003).
- 962 42. Schubbert, S., Shannon, K. & Bollag, G. Hyperactive Ras in developmental disorders and
963 cancer. *Nat. Rev. Cancer* **7**, 295–308 (2007).
- 964 43. Qi, H., Dong, C., Chung, W. K., Wang, K. & Shen, Y. Deep Genetic Connection Between
965 Cancer and Developmental Disorders. *Hum. Mutat.* **37**, 1042–1050 (2016).
- 966 44. Yavuz, B. R. *et al.* Neurodevelopmental disorders and cancer networks share pathways, but
967 differ in mechanisms, signaling strength, and outcome. *Npj Genomic Med.* **8**, 1–14 (2023).
- 968 45. Amberger, J. S., Bocchini, C. A., Scott, A. F. & Hamosh, A. OMIM.org: leveraging
969 knowledge across phenotype-gene relationships. *Nucleic Acids Res.* **47**, D1038–D1043 (2019).
- 970 46. Deng, Z. *et al.* Familial Primary Pulmonary Hypertension (Gene PPH1) Is Caused by
971 Mutations in the Bone Morphogenetic Protein Receptor–II Gene. *Am. J. Hum. Genet.* **67**, 737–744
972 (2000).
- 973 47. Momose, Y. *et al.* De novo mutations in the *BMPR2* gene in patients with heritable
974 pulmonary arterial hypertension. *Ann. Hum. Genet.* **79**, 85–91 (2015).
- 975 48. Evans, J. D. W. *et al.* *BMPR2* mutations and survival in pulmonary arterial hypertension: an
976 individual participant data meta-analysis. *Lancet Respir. Med.* **4**, 129–137 (2016).
- 977 49. Bashamboo, A., Bignon-Topalovic, J., Moussi, N., McElreavey, K. & Brauner, R. Mutations
978 in the Human *ROBO1* Gene in Pituitary Stalk Interruption Syndrome. *J. Clin. Endocrinol. Metab.*
979 **102**, 2401–2406 (2017).
- 980 50. Liu, Z. & Chen, X. A Novel Missense Mutation in Human Receptor Roundabout-1 (*ROBO1*)
981 Gene Associated with Pituitary Stalk Interruption Syndrome. *J. Clin. Res. Pediatr. Endocrinol.* **12**,
982 212–217 (2020).
- 983 51. Liu, C. *et al.* Deficiency of primate-specific *SSX1* induced asthenoteratozoospermia in
984 infertile men and cynomolgus monkey and tree shrew models. *Am. J. Hum. Genet.* **110**, 516–530
985 (2023).

- 986 52. Alankarage, D. *et al.* Myhre syndrome is caused by dominant-negative dysregulation of
987 SMAD4 and other co-factors. *Differentiation* **128**, 1–12 (2022).
- 988 53. Rentzsch, P., Schubach, M., Shendure, J. & Kircher, M. CADD-Splice—improving genome-
989 wide variant effect prediction using deep learning-derived splice scores. *Genome Med.* **13**, 31
990 (2021).
- 991 54. Luyckx, I., Verstraeten, A., Goumans, M.-J. & Loeys, B. SMAD6-deficiency in human
992 genetic disorders. *Npj Genomic Med.* **7**, 1–11 (2022).
- 993 55. Sherwood, K. *et al.* Germline de novo mutations in families with Mendelian cancer
994 syndromes caused by defects in DNA repair. *Nat. Commun.* **14**, 3636 (2023).
- 995 56. Garcia-Salinas, O. I. *et al.* The impact of ancestral, environmental and genetic influences on
996 germline de novo mutation rates and spectra. 2024.05.17.594464 Preprint at
997 <https://doi.org/10.1101/2024.05.17.594464> (2024).
- 998 57. Balmain, A. The critical roles of somatic mutations and environmental tumor-promoting
999 agents in cancer risk. *Nat. Genet.* **52**, 1139–1143 (2020).
- 1000 58. Turner, T. N. *et al.* denovo-db: a compendium of human de novo variants. *Nucleic Acids Res.*
1001 **45**, D804–D811 (2017).
- 1002 59. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in
1003 141,456 humans. *Nature* **581**, 434–443 (2020).
- 1004 60. Kong, A. *et al.* Rate of de novo mutations, father’s age, and disease risk. *Nature* **488**, 471–
1005 475 (2012).
- 1006 61. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat. Genet.* **48**, 126–
1007 133 (2016).
- 1008 62. Jónsson, H. *et al.* Parental influence on human germline de novo mutations in 1,548 trios
1009 from Iceland. *Nature* **549**, 519–522 (2017).
- 1010 63. Kaplanis, J. *et al.* Genetic and chemotherapeutic influences on germline hypermutation.
1011 *Nature* **605**, 503–508 (2022).
- 1012 64. Kingdom, R. *et al.* Rare genetic variants in genes and loci linked to dominant monogenic
1013 developmental disorders cause milder related phenotypes in the general population. *Am. J. Hum.*

- 1014 *Genet.* **109**, 1308–1316 (2022).
- 1015 65. Giannoulatou, E. *et al.* Whole-genome sequencing of spermatocytic tumors provides insights
1016 into the mutational processes operating in the male germline. *PLOS ONE* **12**, e0178169 (2017).
- 1017 66. Armitage, P. & Doll, R. The age distribution of cancer and a multi-stage theory of
1018 carcinogenesis. *Br. J. Cancer* **8**, 1–12 (1954).
- 1019 67. Nowell, P. C. The clonal evolution of tumor cell populations. *Science* **194**, 23–28 (1976).
- 1020 68. Cimadomo, D. *et al.* Impact of Maternal Age on Oocyte and Embryo Competence. *Front.*
1021 *Endocrinol.* **9**, 327 (2018).
- 1022 69. Paul, C. & Robaire, B. Ageing of the male germ line. *Nat. Rev. Urol.* **10**, 227–234 (2013).
- 1023 70. Moore, J. E. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse
1024 genomes. *Nature* **583**, 699–710 (2020).
- 1025 71. Gao, Y.-B. *et al.* Genetic landscape of esophageal squamous cell carcinoma. *Nat. Genet.* **46**,
1026 1097–1102 (2014).
- 1027 72. Waddell, N. *et al.* Whole genomes redefine the mutational landscape of pancreatic cancer.
1028 *Nature* **518**, 495–501 (2015).
- 1029 73. Aaltonen, L. A. *et al.* Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- 1030 74. Nguyen, B. *et al.* Genomic characterization of metastatic patterns from prospective clinical
1031 sequencing of 25,000 patients. *Cell* **185**, 563-575.e11 (2022).
- 1032 75. Kakiuchi, N. & Ogawa, S. Clonal expansion in non-cancer tissues. *Nat. Rev. Cancer* **21**, 239–
1033 256 (2021).
- 1034 76. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
1035 Preprint at <https://doi.org/10.48550/arXiv.1303.3997> (2013).
- 1036 77. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**,
1037 308–311 (2001).
- 1038 78. Zhang, F. *et al.* Ancestry-agnostic estimation of DNA sample contamination from sequence
1039 reads. *Genome Res.* **30**, 185–194 (2020).
- 1040 79. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008 (2021).
- 1041 80. Subramanian, S. & Kumar, S. Neutral Substitutions Occur at a Faster Rate in Exons Than in

- 1042 Noncoding DNA in Primate Genomes. *Genome Res.* **13**, 838–844 (2003).
- 1043 81. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models Using
1044 lme4. *J. Stat. Softw.* **67**, 1–48 (2015).
- 1045 82. Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat. Genet.*
1046 **47**, 1402–1407 (2015).
- 1047 83. Martinelli, S. *et al.* Molecular Diversity and Associated Phenotypic Spectrum of Germline
1048 CBL Mutations. *Hum. Mutat.* **36**, 787–796 (2015).
- 1049 84. Khadka, P. *et al.* PPM1D mutations are oncogenic drivers of de novo diffuse midline glioma
1050 formation. *Nat. Commun.* **13**, 604 (2022).
- 1051 85. Frankish, A. *et al.* GENCODE 2021. *Nucleic Acids Res.* **49**, D916–D923 (2021).
- 1052 86. The UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic*
1053 *Acids Res.* **51**, D523–D531 (2023).
- 1054 87. Ou, J. & Zhu, L. J. trackViewer: a Bioconductor package for interactive and integrative
1055 visualization of multi-omics data. *Nat. Methods* **16**, 453–454 (2019).
- 1056 88. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).