1  **Genetic and chemotherapeutic causes of germline hypermutation**
2

3  **Authors:** Joanna Kaplanis[1], Benjamin Ide[2],  Rashesh Sanghvi[1], Matthew Neville[1], Petr
4  Danecek[1], Tim Coorens[1], Elena Prigmore[1],  Patrick Short[1], Giuseppe Gallone[1], Jeremy
5  McRae[1], Chris Odhams[3], Loukas Moutsianas[3], Genomics England Research Consortium,
6  Jenny Carmichael[4], Angela Barnicoat[5], Helen Firth[1,4], Patrick O'Brien[2], Raheleh Rahbari[1],
7  Matthew Hurles[1]
8  **Affiliations:**
9  [1] Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK
10  [2] Department of Biological Chemistry, University of Michigan, Ann Arbor, Michigan USA
11  [3] Genomics England, London, UK
12  [4] East Anglian Medical Genetics Service, Cambridge University Hospitals, Cambridge, UK
13  [5] North East Thames Regional Genetics Service, Great Ormond Street Hospital, London, UK
14

15  **Summary:**
16  Mutation in the germline is the source of all evolutionary genetic variation and a cause of
17  genetic disease. Previous studies have shown parental age to be the primary determinant of
18  the number of new germline mutations seen in an individual's genome. Here we analysed
19  the genome-wide sequences of 21,879 families with rare genetic diseases and identified 12
20  hypermutated individuals with between two and seven times more *de novo* single nucleotide
21  variants (dnSNVs) than expected. In most of these families (8/12) the excess mutations
22  could be attributed to the father. We determined that two of these families had genetic
23  drivers of germline hypermutation, with the fathers carrying damaging genetic variation in
24  known DNA repair genes, causing distinctive mutational signatures. For five families, by
25  analysing clinical records and mutational signatures, we determined that paternal exposure
26  to chemotherapeutic agents prior to conception was a key driver of hypermutation. Our
27  results suggest that the germline is well protected from mutagenic effects, hypermutation is
28  rare and relatively modest in degree and that most hypermutated individuals will not have a
29  genetic disease.

30  # Introduction

31  Germline mutagenesis is the source of all genetic variation which drives evolution and
32  generates disease-causing variants. The average number of *de novo* mutations (DNMs)
33  generating single nucleotide variants (SNVs) is estimated to be 60-70 per human genome
34  per generation, but little is known about germline hypermutated individuals with unusually
35  large numbers of DNMs [1–3]. The human germline mutation rate is not a constant, but varies
36  between individuals, families and populations and has evolved over time just like any other
37  phenotype [45–8]. Parental age explains a large proportion of variance for single nucleotide
38  variants (SNVs), indels and short tandem repeats (STRs) [3,9,10] It has been estimated that
39  there is an increase of ~2 DNMs for every additional year in father's age and a more subtle
40  increase of ~0.5 DNMs for every additional year in mother's age [3,11]. Subtle differences have
41  also been observed between the maternal and paternal mutational spectra and may be
42  indicative of different mutagenic processes [12–15]. Different mutational mechanisms can leave
43  distinct mutational patterns. These combinations of mutation types can be decomposed from
44  mutational spectra into 'mutational signatures' [16,17]. There are currently >100 somatic
45  mutational signatures that have been identified across a wide variety of cancers of which half
46  have been attributed to endogenous mutagenic processes or specific mutagens [18,19]. The
47  majority of germline mutation can be explained by two of these signatures, termed signature

48  1 (SBS1) , likely due to deamination of 5-methylcytosine [20], and signature 5 (SBS5), thought
49  to be a pervasive and relatively clock-like endogenous process. Both signatures are
50  ubiquitous among normal and cancer cell types[21,22] and have been reported previously in
51  trio-studies[13]. The impact of environmental mutagens has been well established in the soma
52  but is not as well understood in the germline[23,24]. Environmental exposures in parents, such
53  as ionising radiation, can influence the number of mutations transmitted to offspring[25–27].
54  Individual mutation rates can also be influenced by genetic background. With regards to
55  somatic mutation, thousands of inherited germline variants have been shown to increase
56  cancer risk[28–30]. Many of these variants are in genes encoding components of DNA repair
57  pathways which, when impaired, lead to an increased number of somatic mutations.
58  However it is not known whether variants in known somatic mutator genes can influence
59  germline mutation rates. There are a handful of examples where genetic background has
60  been shown to impact the germline mutation rate of STRs, minisatellites and translocations,
61  often in cis, rather than genome-wide [31–3435].

62      An elevated germline mutation rate can have a significant impact on the health of
63  subsequent generations. Increasing germline mutation rate results in an increased risk of
64  offspring being born with a genetic disorder caused by a DNM[36]. Long-term effects of
65  mutation rate differences as a result of mutation accumulation have been demonstrated in
66  mice to have effects on reproduction and survival rates and there may be a similar impact in
67  humans [37,38].

68      While we have started to explain the general properties of germline mutations, little is
69  known about rare outliers with extreme mutation rates. *De novo* mutations are a substantial
70  cause of rare genetic disorders and cohorts of patients with such disorders are enriched for
71  DNMs overall and are more likely to include germline hypermutated individuals[11,39]. To this
72  end we sought to identify germline hypermutated individuals in ~20,000 sequenced parent
73  offspring trios from two rare disease cohorts. We identified genetic or environmental causes
74  of this hypermutation and estimated how much variation in germline mutation rate this may
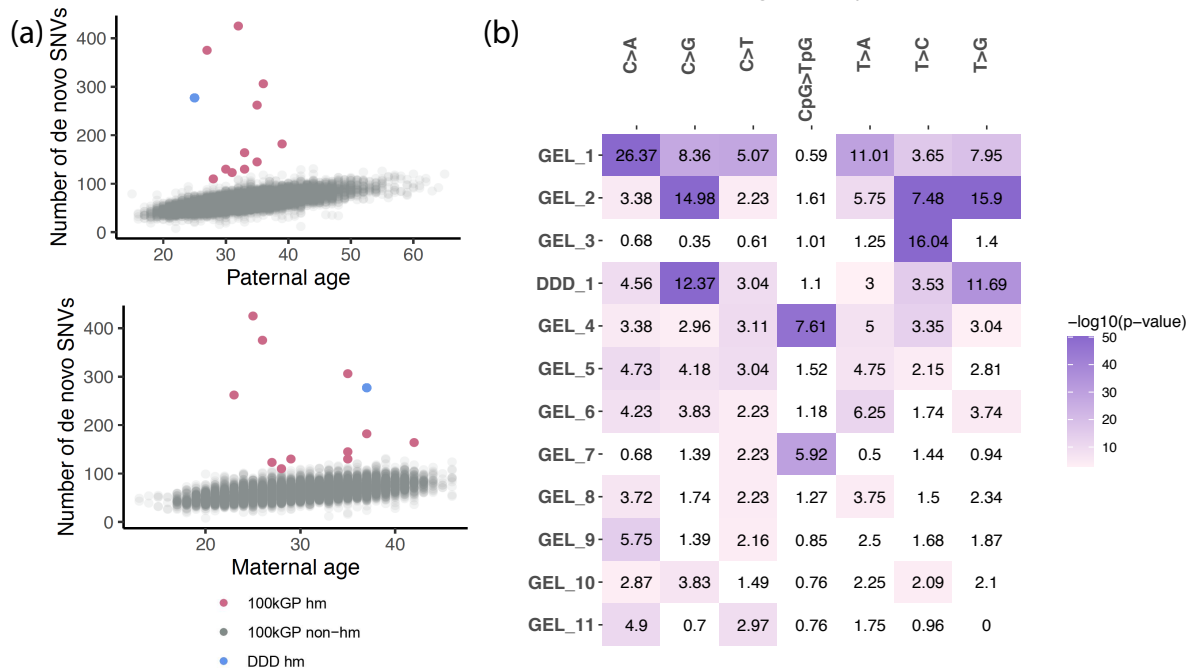75  explain.

# Results

**Identifying germline hypermutated individuals in rare disease cohorts**
78  We sought to identify germline hypermutated individuals in two separate cohorts: 7,930
79  exome-sequenced parent offspring trios from the Deciphering Developmental Disorders
80  (DDD) Study and 13,949 whole-genome sequenced parent offspring trios in the rare disease
81  arm of the 100,000 Genome Project (100kGP). We selected nine trios from the DDD study
82  with  the largest number of exonic DNMs in the offspring, given their parental ages, which
83  were subsequently whole genome sequenced at >30X coverage to characterise DNMs
84  genome-wide. In the 100kGP cohort, we performed extensive filtering of the DNMs which
85  resulted in a total of 903,525 *de novo* SNVs (dnSNVs) and 72,110 *de novo* indels (dnIndels).
86  The median number of DNMs per individual was 62 for dnSNVs and 5 for dnIndels (median
87  paternal and maternal ages of 33 and 30) (Supplemental Figure 1).

88      Parental age explains the majority of variance in numbers of germline mutations
89  observed in offspring and is important to control for when examining additional sources of
90  variation [3]. We observed an increase in total number of dnSNVs of 1.28 dnSNVs/year of
91  paternal age (CI:1.24-1.32, p<$10^{-300}$, Negative binomial regression) and an increase of 0.35
92  dnSNVs/year of maternal age (CI: 0.30-0.39, p = 3.0 $\times 10^{-49}$, Negative Binomial regression)
93  (Figure 1a). We were able to phase 241,063 dnSNVs and found that 77% of phased DNMs
94  were paternal in origin, which agrees with previous estimates[12–14]. Estimates of the parental

95 age effect in the phased mutations were not significantly different to the unphased results:
96 1.23 paternal dnSNVs/year of paternal age (CI: 1.14-1.32, p =$1.6\times10^{-158}$ ) and 0.38 maternal
97 dnSNVs/year of maternal age ( CI: 0.35,0.41, p = $6.6\times10^{-120}$) (Supplemental Figure 2b).
98 Paternal and maternal age were also significantly associated with the number of dnIndels:
99 an increase of 0.071 dnIndels/year of paternal age (CI: 0.062-0.080, p = $8.3\times10^{-56}$,
100 Supplemental Figure 2a) and a smaller increase of 0.019 dnIndels/year of maternal age (CI:
101 0.0085-0.029 p = $3.4 \times10^{-4}$, Supplemental Figure 2a). The ratio of paternal to maternal
102 mutation increases for SNVs and indels were very similar, 3.7 for SNVs and 3.8 for indels.
103 The proportion of *de novo* mutations that phased paternally increased by 0.0017 for every
104 year of paternal age (p = $3.37 \times10^{-38}$, Binomial regression, Supplemental Figure 3).
105 However, the effect size is small and the proportion of DNMs that phase paternally in the
106 youngest fathers is ~0.75 and so the paternal age effect alone does not fully explain the
107 strong paternal bias [14]. We compared the mutational spectra of the phased DNMs and found
108 that maternally derived DNMs have a significantly higher proportion of C>T mutations (0.27
109 maternal vs 0.22 paternal, p = $3.24\times10^{-80}$, Binomial test) , while paternally derived DNMs
110 have a significantly higher proportion of C>A, T>G and T>C mutations (C>A: 0.08 maternal
111 vs 0.10 paternal, p = $4.6\times10^{-23}$; T>G 0.06 vs 0.7, p = $6.8\times10^{-28}$; T>C 0.25 vs 0.26, p =
112 $1.6\times10^{-5}$; Binomial test, Supplemental Figure 4a). These mostly agree with previous studies
113 although the difference in T>C mutations was not previously significant[12]. The majority of
114 both paternal and maternal mutations could be explained by Signature 1 and 5, with a
115 slightly higher contribution of signature 1 in paternal mutations (0.16 paternal vs 0.15
116 maternal, chi-squared test p = $2.0 \times 10^{-5}$,, Supplemental Figure 4b).
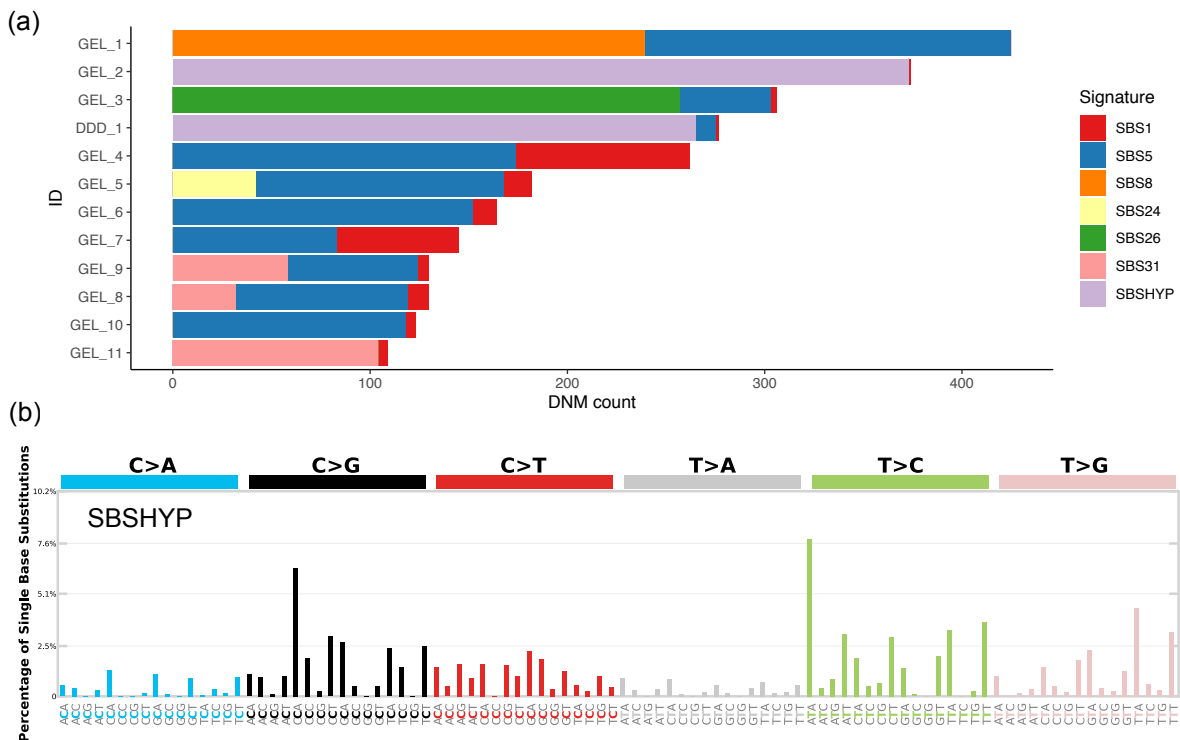


117

118 **Figure 1: Identification of germline hypermutated individual**s (a) Paternal and maternal
119 age vs number of dnSNVs, 100kGP hypermutated individuals are highlighted in pink and
120 DDD hypermutated individual is highlighted in blue (b) Enrichment (observed/expected) of
121 mutation type for hypermutated individuals. Sample names on the y-axis, mutation type on
122 the x-axis. The enrichment is colored by the -log10(enrichment p-value) which was
123 calculated using a Poisson test comparing the average number of mutations in each type
124 across all individuals in the 100kGP cohort. White coloring indicates no statistically
125 significant enrichment (p-value <0.05/12*7).

| ID | Number of dnSNVs/ dnIndels | Child age | Paternal age | Maternal age | SNV p-value | Indel p-value | TS bias | Phase (P,M) | Potential source of hypermutation |
|---|---|---|---|---|---|---|---|---|---|
| GEL_1 | 425/16 | 5-10 | 30-35 | 20-25 | 4.2e-90 | 9.4e-05 | 2.1e-40 | 129,1*** | Paternal DNA repair defect; homozygous stop-gain *XPC* variant |
| GEL_2 | 375/5 | 10-15 | 25-30 | 25-30 | 2.3e-83 | 0.43 | 0.22 | 106,7*** | Paternal chemotherapy; Nephrotic syndrome: Cyclophosphamide, Chlorambucil |
| GEL_3 | 306/4 | 0-5 | 35-40 | 30-35 | 2.5e-44 | 0.73 | 0.86 | 87,5*** | Paternal DNA repair defect; homozygous missense *MPG* variant |
| DDD_1 | 277/6 | 6 | 25 | 37 | NA | NA | 3.3e-03 | 72,4*** | Paternal chemotherapy; Hodgkins Lymphoma: ABVD, IVE |
| GEL_4 | 262/12 | 10-15 | 30-35 | 20-25 | 1.7e-37 | 0.007 | 0.070 | 36,32*** | Post-zygotic hypermutation |
| GEL_5 | 182/8 | 0-5 | 35-40 | 35-40 | 8.4e-14 | 0.19 | 0.15 | 63,4*** | Paternal chemotherapy; SLE: drugs unknown |
| GEL_6 | 164/7 | 0-5 | 30-35 | 40-45 | 9.8e-13 | 0.25 | 0.066 | 38,3* | Unknown |
| GEL_7 | 145/9 | 0-5 | 30-35 | 30-35 | 2.4e-09 | 0.08 | 0.02 | 24,16* | Post-zygotic hypermutation |
| GEL_8 | 130/6 | 20-25 | 25-30 | 25-30 | 2.1e-09 | 0.31 | 1.00 | 31,11 | Paternal chemotherapy; Testicular cancer: drugs unknown |
| GEL_9 | 130/5 | 5-10 | 30-35 | 30-35 | 1.2e-07 | 0.53 | 0.016 | 46,2** | Paternal chemotherapy; Testicular cancer: BEP |
| GEL_10 | 123/5 | 10-15 | 30-35 | 25-30 | 5.3e-08 | 0.48 | 0.082 | 38,0*** | Unknown |
| GEL_11 | 110/5 | 10-15 | 25-30 | 25-30 | 8.2e-07 | 0.44 | 6.9e-06 | 28,1* | Paternal chemotherapy; Cancer of long bones, intestinal tract, lung (secondary): Drugs unknown |

126 **Table 1: Properties and possible hypermutation sources for 12 germline**
127 **hypermutated individuals.** Eleven individuals were identified in 100kGP as having a
128 significantly large number of dnSNVs (GEL_1-GEL_11) and one individual identified in the
129 DDD study (DDD_1). The DNM counts are for autosomal DNMs. Child age refers to age
130 when sample was taken. Paternal and maternal age refer to age at child's birth. All ages are
131 given as 5 year ranges for 100kGP individuals and the exact age for the DDD individuals.
132 SNV and indel p-value is from testing the number of dnSNVs and dnIndels compared to
133 what we would expect after accounting for parental age. TS bias: transcriptional strand bias
134 p-value for dnSNVS. Phase (P,M): the number of dnSNVs that phase paternally (P) and
135 maternally (M) with significance indicator for how different this ratio is compared to the
136 observed proportion across all DNMs that phase paternally in 100kGP (0.77) using a
137 Binomial test (*p<0.1, **p<0.01,***p<0.001). We have detailed the parental cancer and
138 chemotherapy drugs received when relevant. Treatments abbreviations: BEP (Bleomycin,
139 etoposide and platinum), ABVD (Bleomycin-Dacarbazine-Doxorubicin-Vinblastine) and IVE
140 (Iphosphamide, epirubicin and etoposide).
141

142     We identified 12 germline hypermutated individuals after accounting for parental age
143  (see Methods): 11 from the 100kGP cohort and 1 from the DDD cohort (Figure 1a, Table 1).
144  The number of DNMs genome-wide for each of the 12 hypermutated individuals ranged from
145  110-425 dnSNVs, which corresponds to a fold increase of 1.7-6.5 compared to the median
146  number of dnSNVs per individual across the 100kGP cohort. Two of these individuals also
147  had a significantly increased number of dnIndels (Table 1). The mutational spectra across
148  these hypermutated individuals varied dramatically (Figure 1b, Supplemental Figure 5,
149  Supplemental Table 1) and after extracting mutational signatures we found that while many
150  of the mutations mapped onto several known somatic signatures (from the Catalogue of
151  Somatic Mutations in Cancer (COSMIC)[40]), a novel mutational signature, termed SBSHYP,
152  was also extracted (Figure 2a,b, Supplemental Table 2). In addition to mutational spectra,
153  we evaluated the parental phase, transcriptional strand bias (Supplemental Figure 6) and the
154  distribution of the variant allele fraction (VAF) for these mutations (Supplemental Figure 7).
155  Upon examining these properties, we identified three potential sources of germline
156  hypermutation: paternal defects in DNA repair genes, paternal exposure to
157  chemotherapeutics and post-zygotic mutational factors.



158
159  **Figure 2: Mutational signatures in 12 germline hypermutated individuals** (a)
160  Contributions of mutational signatures extracted using SigProfiler and decomposed on to
161  known somatic mutational signatures as well as the novel signature (SBSHYP) identified in
162  both DDD_1 and GEL_2. Summary of signatures: SBS1 and SBS5 are known germline
163  signatures; SBS8 associated with TC-NER; SBS26 associated with defective MMR; SBS31
164  associated with platinum drug treatment; SBS24 is associated with aflatoxin exposure.  (b)
165  Trinucleotide context mutational profile of novel extracted mutational signature SBSHYP
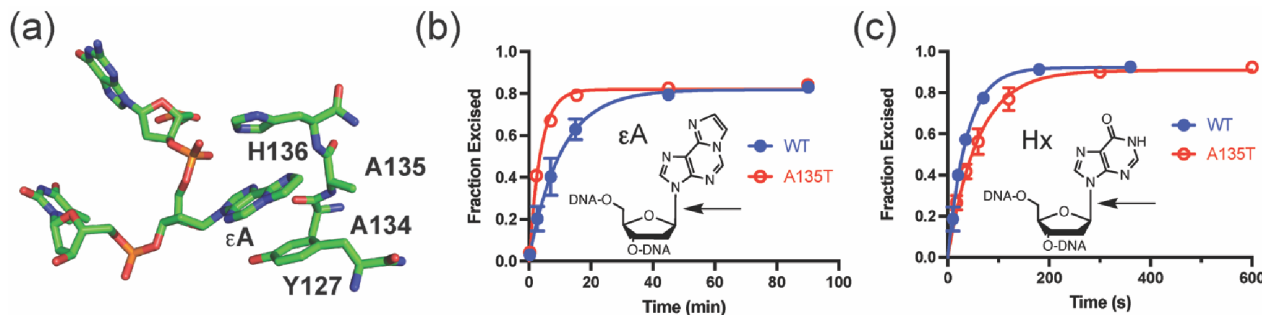166
167  **Paternal defects in DNA repair**
168  For eight of the twelve individuals, the DNMs phased paternally significantly more than
169  expected given the overall ratio of paternal:maternal mutation in the 100kGP cohort
170  (p<0.05/12, Binomial test, Table 1). This implicates the paternal germline as the source of
171  the hypermutation. Two of these fathers carry rare homozygous nonsynonymous variants in

172  known DNA repair genes (Supplemental Table 3). Defects in DNA repair are known to
173  increase the mutation rate in the soma and may have a similar effect in the germline. GEL_1
174  has the largest number of DNMs of all individuals, a 6.5-fold enrichment, and a significantly
175  increased number of dnIndels. The mutational spectra demonstrates a high enrichment of
176  C>A and T>A mutations (Figure 1b) and we observed a large contribution from Somatic
177  Mutational Signature 8 (Figure 2a). This signature is associated with transcription-coupled
178  nucleotide excision repair (TC-NER) and typically presents with transcriptional strand bias.
179  This agrees with the strong transcriptional strand bias observed in GEL_1 (p = 2.1 $\times 10^{-40}$,
180  Poisson test, Supplemental Figure 6). The father has a rare homozygous nonsense variant
181  in the gene *XPC* (Table 1, Supplemental Table 3) which is involved in the early stages of the
182  nucleotide-excision repair (NER) pathway. The paternal variant is annotated as pathogenic
183  for xeroderma pigmentosum in ClinVar and clinical follow-up confirmed that the father had
184  already been diagnosed with this disorder. Patients with xeroderma pigmentosum have a
185  high risk of developing skin cancer due to their impaired ability to repair UV damage and are
186  also known to be at a higher risk of developing other cancers [41,42]. *XPC* deficiency has been
187  associated with a similar mutational spectrum to the one we observe in GEL_1[43] and xpc
188  deficiency in mice has been shown to increase the germline mutation rate at two STR loci [44].
189      GEL_3 has a ~5-fold enrichment of the number of dnSNVs. These dnSNVs exhibit a
190  very distinct mutational spectrum with a ~17-fold increase in T>C mutations but no
191  significant enrichment for any other mutation type (Figure 2b, Supplemental Figure 5d).
192  Extraction of mutational signatures revealed that the majority of mutations mapped onto
193  Somatic Mutational Signature 26 which has been associated with defective mismatch repair.
194  The father has a rare homozygous missense variant in the gene *MPG* (Table 1,
195  Supplemental Table 4). *MPG* encodes N-methylpurine DNA glycosylase (also known as
196  alkyladenine-DNA glycosylase - AAG) which is involved in the recognition of base lesions,
197  including alkylated and deaminated purines, and initiation of the base-excision repair (BER)
198  pathway. The *MPG* variant is rare in gnomAD  (allele frequency=9.8 $\times 10^{-5}$ , no observed
199  homozygotes) and is predicted to be pathogenic by the Combined Annotation Dependent
200  Depletion (CADD) score (CADD score = 27.9) and the amino acid residue is fully conserved
201  across 172 aligned protein sequences from VarSite [45,46]. In the context of the protein, the
202  variant amino-acid forms part of the substrate binding pocket and likely affects substrate
203  specificity (Figure 3a). *MPG* has not yet been described as a cancer susceptibility gene, but
204  studies in yeast and mice have demonstrated variants in this gene, and specifically the
205  substrate binding pocket, can lead to a mutator phenotype [47,48]. We explored the functional
206  impact of the observed A135T variant using *in vitro* assays (Methods, Supplemental Figures
207  8, 9). The A135T variant caused a two-fold decrease in excision efficiency of the
208  deamination product hypoxanthine (Hx) in both the T and C contexts (Figure 3c,
209  Supplemental Figure 9), with a small increase in excision efficiency of an alkylated adduct
210  1,N6-ethenoadenine ($\varepsilon$A) in both the T and C contexts (Figure 3b, Supplemental Figure 9).
211  The maximal rate of excision is increased by 2-fold for $\varepsilon$A which is among the largest
212  increases that have been observed for 15 reported MPG variants (Supplemental Table 4).
213  Another variant, N169S, which also shows an increase in N-glycosidic bond cleavage with
214  the $\varepsilon$A substrate has been established as a mutator in yeast[48,49]. These assays confirm that
215  the A135T substitution alters the MPG binding pocket and changes the activity towards
216  different DNA adducts. MPG acts on a wide variety of DNA adducts and further functional
217  characterisation and mechanistic studies are required to link the observed T>C germline
218  mutational signature to the aberrant processing of a specific class of DNA adducts.
219      GEL_6 has 164 dnSNVs and has a larger contribution of paternally phased mutations
220  than expected (38:3 paternal:maternal, p = 0.022, Binomial test) however the father does not

221 have any nonsynonymous variants in DNA repair genes nor has undergone any
222 chemotherapeutic treatment. This unexplained cause of hypermutation could be due to a
223 paternal variant in a gene currently not associated with DNA repair, a paternal mutator
224 variant that is only present in the germline and not the blood, or a rare gene-by-environment
225 interaction.
226



227
228 **Figure 3. A135T substitution alters the DNA glycosylase activity of MPG.** (a) Active site
229 view of MPG bound to εA-DNA from pdb 1EWN. A135 and H136 form the binding pocket for
230 the flipped-out base lesion, which is bracketed by Y127 on the opposing face. (b) Single
231 turnover excision of εA from εA•T is 2 fold faster for A135T (red) than for WT (blue) MPG. (c)
232 Single turnover excision of Hx from Hx•T is slower for A135T (red) as compared to WT (blue)
233 MPG. Arrows indicate the N-glycosidic bond that is cleaved by MPG. Each data point is the
234 mean ± SD (N≥3) (see Supplemental Figure 9 for complete kinetic analysis).
235

**Parental treatment with chemotherapy prior to conception**
237 Three hypermutated individuals (GEL_8, GEL_9 and GEL_11) have a contribution from
238 somatic mutational signature 31 (SBS31) (Figure 2a), which has been associated with
239 treatment with platinum-based drugs such as cisplatin[16]. The phased dnSNVs in GEL_9 and
240 GEL_11 are paternally biased (46 paternal: 2 maternal, p = 0.0014; 28 paternal:1 maternal,
241 p = 0.012; Binomial test, Table 1), and the dnSNVs in GEL_11, which has the largest
242 contribution of SBS31, exhibit a significant transcriptional strand bias, as expected for this
243 signature (p = 6.9 ×10^{-6}, Table 1, Supplemental Figure 6). All three fathers have a history
244 of cancer and chemotherapeutic treatment prior to conception as recorded in their available
245 hospital episode records. The father of GEL_11 was diagnosed and received
246 chemotherapeutic treatment for osteosarcoma, lung cancer, and cancer of the intestinal tract
247 within 5 years prior to conception. Cisplatin is a commonly used chemotherapeutic for
248 osteosarcoma and lung cancer. Platinum-based drugs damage DNA by causing covalent
249 adducts. Cisplatin mainly reacts with purine bases, forming intrastrand crosslinks which can
250 be repaired by NER or bypassed by translesion synthesis which may, in turn, induce single
251 base substitutions[50]. The fathers of GEL_8 and GEL_9 both have a history of testicular
252 cancer where cisplatin is the most commonly administered chemotherapeutic.
253 GEL_2 and DDD_1 have a similar number of dnSNVs, which are significantly
254 paternally biased (Table 1). The mutational spectra of the DNMs in these individuals are very
255 similar and share a novel mutational signature (SBSHYP) that is characterised by an
256 enrichment of C>G and T>G mutations (Figure 2a,b) and does not map on to any previously
257 described signature observed in somatic mutations (as described in COSMIC) or in
258 response to mutagenic exposure[24,51,52] (Supplemental Figure 10a). The fathers of these
259 individuals do not have putative damaging variants in any DNA repair genes and do not
260 share rare nonsynonymous variants in any other gene. Both fathers received
261 chemotherapeutic treatment prior to conception including nitrogen mustard alkylating agents
262 (Supplemental Table 4), although with different members of this class of chemotherapies,

263  therefore we strongly suspect this class of chemotherapies to be the cause of this novel
264  mutational signature. Experimental studies of a subset of alkylating agents have shown them
265  to have diverse mutational signatures[24,51,52] (Supplemental Figure 10b).
266      GEL_5 has 182 dnSNVs and a significant paternal bias in the phased dnSNVs (p =
267  $5.8 \times 10^{-4}$, Binomial Test, Table 1). The father of GEL_5 has a diagnosis of Systemic Lupus
268  Erythematosus (SLE) and received a course of chemotherapy nine years prior to the
269  conception of the child however the dnSNVs do not map onto any known chemotherapeutic
270  mutational signatures (Figure 1b, Figure 2a). There is a contribution of SBS24 which is
271  associated with aflatoxin exposure in cancer blood samples, however there is no evidence of
272  exposure in the father's hospital records[22]
273      We assessed how parental cancer and exposure to chemotherapy might impact
274  germline mutation rate more generally by systematically examining hospital episode
275  statistics across the 100kGP cohort for ICD10 codes related to cancer and chemotherapy
276  that were recorded prior to the conception of the child. We identified 27 fathers (0.9%) who
277  had a history of cancer, 7 of which had testicular cancer (Supplemental Table 6). The
278  offspring of these 27 fathers did not have a significantly increased number of dnSNVs after
279  correcting for parental age (p = 0.73, Wilcox test). This is a small number of fathers so this is
280  not well powered and it is not known definitively how many of these fathers were treated with
281  chemotherapy (6 of the 27 had chemotherapy-related ICD10 codes). Treatment exposure
282  may predate the availability of digitised hospital records and there was limited information on
283  whether conception may have been achieved using sperm stored prior to treatment. While
284  the total number of dnSNVs across all the children is not significantly increased, two of the
285  27 fathers had hypermutated children which is a significant enrichment compared to those
286  fathers who do not have a recorded history of cancer (2/27 vs 9/2891, p = 0.0043, Fisher
287  exact test). This is likely a conservative p-value as we know that two other hypermutated
288  individuals have fathers who have been treated with chemotherapy however did not fall into
289  this group due to the filtering criteria as we only considered fathers who had at least one ICD
290  10 code recorded prior to the child's conception (see Methods). A possible confounder could
291  be that fathers who had cancer prior to conception are older and may be more likely to have
292  been exposed to other germline mutagens however these two groups had the same median
293  paternal age (p = 0.77, Wilcoxon test).
294      We performed the same analysis across 5,508 mothers in the 100kGP cohort who
295  had hospital episode records entered prior to the conception of their child and identified 27
296  mothers (0.5%) who had a history of cancer, 9 of whom also had recorded chemotherapy
297  codes. Children whose mothers had a history of cancer had a nominally significant increase
298  of dnSNVs after correcting for parental age and data quality (p = 0.03, Wilcox Test). Mothers
299  who had been diagnosed with cancer were significantly older at the birth of the child
300  compared to those who were not (p = 0.003, Wilcoxon test). Matching on parental age,
301  mothers who had a cancer diagnosis prior to conception had a median increase of 9
302  dnSNVs. Overall there is not an excess of maternally phased DNMs across these individuals
303  (p = 0.44, Binomial test) however there is one individual with nominal significance
304  (MatCancer_23, 22 paternal:14 maternal, p =0.02, Binomial test, Supplemental Table 5).
305      We extracted mutational signatures for all these offspring with a maternal or paternal
306  history of cancer that were not hypermutated and found that only one individual had unusual
307  mutational signatures (Supplemental Figure 11).This individual (PatCancer_10)  has a
308  contribution of mutational signature SBS31 which is associated with treatment with platinum-
309  based drugs (Supplemental Figure 11). Their father was treated for testicular cancer prior to
310  conception and the child has 94 dnSNVs  (p = 0.005, SNV p-value after correcting for

311    parental age) of which 89% phased paternally (p = 0.12, Binomial test, Supplemental Figure
312    11, Supplemental Table 6).
313
314    **Post-zygotic hypermutation**
315    The two hypermutated individuals, GEL_4 and GEL_7, have a ~4 fold and ~2 fold increase
316    in dnSNVs respectively that phase equally between the maternal and paternal
317    chromosomes. The allele balance of the dnSNVs in these individuals was shifted below 0.5
318    (Supplemental Figure 7). In both individuals, the proportion of DNMs with variant allele
319    fraction (VAF) <0.4 was significantly higher compared to all DNMs across all individuals
320    (GEL_4: p = 3.9 $\times 10^{-59}$, GEL_7: p = 8.3$\times 10^{-4}$, Binomial test). These observations indicate
321    that these mutations most likely occurred post-zygotically and are less likely due to a
322    parental hypermutator. Both individuals have a large contribution of mutations from Somatic
323    Mutational Signature 1 (Figure 2a)[40]. The observations in GEL_4 are likely due to clonal
324    haematopoiesis leading to a large number of somatic mutations in the child's blood. The
325    mutational signature associated with haematopoietic stem cells is similar to SBS1and we
326    identified a mosaic *de novo* missense mutation in the gene *ETV6*. Mutations in *ETV6* are
327    associated with Leukemia and Thrombocytopenia [53]. GEL_4 has several blood related
328    clinical phenotypes such as abnormality of blood and blood-forming tissues and
329    myelodysplasia. We do not observe similar phenotypes in GEL_7, nor did we identify a
330    possible genetic driver of clonal haematopoiesis, and the child was one year old at
331    recruitment. For this individual we considered the possibility that a maternal mutator variant
332    protein may be impacting the mutation rate in the first few cell divisions. We identified a
333    maternal mosaic missense variant in *TP53* which is annotated as pathogenic in ClinVar for
334    Li-Fraumeni syndrome which is characterised by a predisposition to cancer however it is
335    unknown if this variant is also present in the maternal germline and, if present, whether it
336    would be likely to have a germline mutagenic effect [54]. This variant is not observed in the
337    child.
338
339    **Fraction of germline mutation rate variation explained**
340    We investigated the factors influencing the number of dnSNVs per individual in a subset of
341    7,700 100kGP trios filtered more stringently for data quality (Methods). Using a negative
342    binomial model, accounting for the underlying Poisson variation in germline mutation rate,
343    we estimated that parental age accounts for 69.7% and data quality metrics (eg. read depth,
344    proportion of mapped reads) explain 1.3% of the variance. The variance explained by
345    parental age is smaller than a previous estimate of 95% based on a sample of 78 families[3].
346    To assess whether this could be due to uncertainty in the previous estimate, we performed
347    repeated sampling of 78 trios from the 100kGP and refit the model and found that the
348    estimates of the variance explained by parental age can vary dramatically with this smaller
349    sample size (median of 79%, 95% interval [52-100%]) and that 7% of our simulations had an
350    estimate as or more extreme than 95%.
351          We extended this model to account for germline hypermutation by including a
352    variable for the number of excess dnSNVs in the 11 hypermutated individuals in this cohort.
353    We found this explained an additional 7.1% of variance. This leaves 21.9% (19.7% , 23.8%,
354    Bootstrap 95% confidence interval) of variance for numbers of dnSNV per individual
355    unaccounted for. Both mutagenic exposures and genetic variation in DNA repair genes are
356    implicated here as causes of hypermutation, therefore they may also play a more subtle role
357    in the remaining germline mutation rate variation. In addition, polygenic effects and gene by
358    environment interactions may also contribute.

359      To assess whether rare variants in genes known to be involved in DNA repair
360   pathways impact germline mutation rate more generally, we looked across the whole
361   100kGP cohort. We curated three sets of rare nonsynonymous variants that have increasing
362   likelihoods of impacting germline mutation rate: (i) variants in all DNA repair genes (N=186),
363   (ii) variants in genes encoding components of the DNA repair pathways most likely to create
364   SNVs (N=66) and (iii) a subset of these variants that had previously been associated with
365   cancer (see Methods). We focused primarily on the effect of heterozygous variants (MAF<
366   0.001). In the first set of genes we also considered the impact of rare homozygous variants
367   (MAF<0.01) (the counts were too small to assess in the subsequent groups). There was no
368   statistically significant effect in any of these groups of variants after Bonferroni correction
369   (Supplemental Figure 12, Supplemental Table 7). We examined heterozygous protein-
370   truncating variants (PTVs) in the known cancer mutator gene MBD4 which are associated
371   with a three-fold elevated CpG>TpG mutation rate in tumours. We identified and whole-
372   genome sequenced 13 paternal carriers of MBD4 PTVs from the DDD cohort. We found that
373   these individuals did not have a significant increased number of overall DNMs and there was
374   no significant increase in the number of CpG>TpG mutations (p = 0.56, chi-squared test,
375   Supplemental Figure 13). Power modelling suggested there is unlikely to be more than a
376   22% increase in the CpG mutation rate. This further demonstrates that heterozygous PTVs
377   in known somatic mutator genes may not always have a similar effect in the germline.
378      To explore the polygenic contribution to germline mutation rate, we estimated the
379   residual variation in the number of dnSNVs in offspring that was explained by germline
380   variants after correcting for parental age, data quality and hypermutation status. We
381   estimated this separately for fathers and mothers in the 100kGP cohort using GREML-LDMS
382   [55] stratified by minor allele frequency and LD. We found that maternal germline variation
383   (MAF>0.001) did not explain any residual variation ($h^2$ = 0.07, p = 0.21, GCTA reported
384   results, Supplemental Table 8). We found that paternal variation may contribute a substantial
385   fraction of residual variation ($h^2$ = 0.53 [0.20,0.85], p = 0.09) however this is  concentrated
386   exclusively in low frequency variants (0.001<MAF<0.01, $h^2$ = 0.52 [0.01,0.94]) rather than
387   more common variants (MAF> 0.01, $h^2$ = 0.008 [0,0.38], Supplemental Table 7). This will
388   need further investigation with larger sample sizes.

389   # Discussion

390   Germline hypermutation is an uncommon but important phenomenon. We identified 12
391   hypermutated individuals from over 20,000 parent offspring sequenced trios in the DDD and
392   100kGP cohorts with a 2-7 fold increased number of dnSNVs. It is likely that there are
393   additional, currently undetected, germline hypermutated individuals in the DDD cohort. The
394   stringent strategy we adopted to screen this exome-sequenced cohort for potential
395   hypermutated individuals for subsequent confirmation by genome sequencing will have
396   missed some individuals with hypermutation of 2-7 fold.
397      In two of the 12 hypermutated individuals, the excess mutations appeared to have
398   occurred post-zygotically, however for the majority (n=8) of these hypermutated individuals,
399   the excess dnSNVs phased paternally implicating the father as the source of this
400   hypermutation. For five of these fathers, characteristic mutational signatures and clinical
401   records of cancer treatment prior to conception strongly implicated the mutagenic influence
402   of two different classes of chemotherapeutics: platinum-based drugs (3 families) and
403   mustard-derived alkylating agents (2 families). We also identified likely paternal mutator
404   variants in two hypermutated families. These were rare homozygous missense variants in

405  two known DNA repair genes: *XPC* and *MPG*. Functional and clinical data strongly
406  supported the mutagenic nature of these variants.
407      It is well established that defects in DNA repair genes can increase somatic mutation
408  rates and elevate cancer risk [56]. Our findings imply that germline mutation rates can be
409  similarly affected. However, defects in DNA repair pathways do not always behave similarly
410  in the soma and the germline. We interrogated PTVs in an established somatic mutator
411  gene, *MBD4*, and found they did not have a detectable effect in the germline [57]. We also
412  examined the impact of parental rare nonsynonymous variants in DNA repair genes on the
413  number of DNMs in offspring and did not find a significant difference. To detect more subtle
414  effects of these variants other analytical approaches will need to be explored. Paternal
415  variants that have previously been associated with a cancer phenotype were nominally
416  significant but having one of these variants only amounted to an estimated average increase
417  of ~2 DNMs in the child. If only a subset of these variants have an impact in the germline this
418  would dilute the power to detect a mutagenic effect and it is likely that both larger sample
419  sizes and additional variant curation will be needed to investigate this further. There may
420  also be genes and pathways that impact mutation in the germline more than the soma;
421  uncovering the genes and associated variants in these genes will be more challenging.
422      Germline hypermutation accounted for 7% of the variance in germline mutation rate
423  in the 100kGP rare disease cohort. The ascertainment in this cohort for rare disease in the
424  offspring, together with the causal contribution that germline mutation plays in rare diseases,
425  means that germline hypermutated individuals are likely enriched in this cohort relative to the
426  general population.  As a consequence, our estimate of the contribution of germline
427  hypermutation to the variance in numbers of dnSNVs per individual is likely inflated.
428  However, the absolute risk of a germline hypermutator having a child with a genetic disease
429  is still low. The population average risk for having a child with a severe developmental
430  disorder caused by a *de novo* mutation has been estimated to be 1 in 300 births[11] and so a
431  4-fold increase in DNMs in a child would only elevate this absolute risk to just over 1%.
432  Therefore, we anticipate that most germline hypermutated individuals will not have a rare
433  genetic disease, and germline hypermutation will also be observed in healthy population
434  cohorts.
435      The two genetic causes of germline hypermutation that we identified were both
436  recessive in action. Similarly, most DNA repair disorders act recessively in their cellular
437  mutagenic effects. This implies that genetic causes of germline hypermutation are likely to
438  arise at substantially higher frequencies in populations with high rates of parental
439  consanguinity. In such populations, the overall incidence of germline hypermutation may be
440  higher and the proportion of the variance in the number of dnSNVs per offspring accounted
441  for genetic effects will be higher. We anticipate that studies focused on these populations are
442  likely to identify additional mutations that affect germline mutation rate.
443      We found that, among 7,700 100kGP families, parental age only explained ~70% of
444  the variance in numbers of dnSNVs per offspring, which is substantially smaller than a
445  previous estimate of 95% based on a sample of 78 families[3]. Repeated sampling of 78 trios
446  from the 100kGP showed that estimates of the variance explained by parental age can vary
447  dramatically stochastically and we regard our estimate based on two orders of magnitude
448  more trios to be more reliable, although other differences between the studies such as
449  measurement error and criteria for ascertainment of families might be having a subtle
450  influence. The residual ~20% of variation in numbers of germline dnSNVs per individual
451  remains unexplained by parental age, data quality and hypermutation. We found that rare
452  variants in known DNA repair genes are unlikely to account for a large proportion of this
453  unexplained variance. Heritability analyses suggested that polygenic contributions from

454  common variants (MAF>1%) are unlikely to make a substantive contribution to this variance;
455  however, we observed some evidence that the polygenic contribution of intermediate
456  frequency paternal variants (0.001<MAF<0.01) could be more substantial although larger
457  sample sizes are required to confirm this observation.  A limitation to these heritability
458  analyses is that we use DNMs in offspring as a proxy for individual germline mutation rates.
459  Measuring germline mutation rates more directly by, for example, sequencing hundreds of
460  single gametes per individual, should facilitate better powered association studies and
461  heritability analyses.
462          Environmental exposures are also likely to contribute to germline mutation rate
463  variation. We have observed evidence that certain chemotherapeutics can affect germline
464  mutation rate and targeted studies on the germline mutagenic effects of different
465  chemotherapeutics (e.g. in cancer survivor cohorts) will be crucial in understanding this
466  further. We anticipate that these studies will identify considerable heterogeneity in the
467  germline mutagenic effects of different chemotherapeutics, in part due to differences in the
468  pemeability of the blood-testis barrier to different agents[58], as well as variation in the
469  vulnerability to chemotherapeutic germline mutagenesis by sex and age.   As so few
470  individuals are treated for cancer prior to reproduction, chemotherapeutic exposures will not
471  explain a large proportion of the remaining variation in germline mutation rates however
472  chemotherapeutic mutagenesis has important implications for cancer patients who plan to
473  have children, especially in whether they decide to store unexposed gametes for future use
474  of assisted reproductive technologies.
475          Unexplained hypermutation and additional variance in germline mutation rate may be
476  explained by other environmental exposures. A limitation of this study was the lack of data
477  on non-therapeutic environmental exposures. However, and somewhat reassuringly, the
478  relatively tight distribution of DNMs per person in 100kGP suggests that there are unlikely to
479  be common environmental mutagen exposures in the UK (e.g. cigarette smoking) that
480  causes a substantive (e.g, >1.5 times) fold increase in mutation rates and concomitant
481  disease risk. The germline generally appears to be well protected from large increases in
482  mutation rate. However, including a broader spectrum of environmental exposures in future
483  studies would help to identify more subtle effects and may reveal gene-by-environment
484  interactions.
485
486  **Acknowledgments**

516

# Methods

518

**DNM filtering in 100,000 Genomes Project**

We analysed DNMs called in 13,949 parent offspring trios from 12,609 families from the rare disease programme of the 100,000 Genomes Project. The rare disease cohort includes individuals with a wide array of diseases including neurodevelopmental disorders, cardiovascular disorders, renal and urinary tract disorders, ophthalmological disorders, tumour syndromes, ciliopathies and others. These are described in more detail in previous publications [59,60]. The cohort was whole genome sequenced at ~35X coverage and variant calling for these families was performed via the Genomics England rare disease analysis pipeline. The details of sequencing and variant calling have been previously described [60]. DNMs were called by the Genomics England Bioinformatics team using the Platypus variant caller[61]. These were selected to optimise various properties including the number of DNMs per person being approximately what we would expect, the distribution of the VAF of the DNMs to be centered around 0.5 and the true positive rate of DNMs to be sufficiently high as calculated from examining IGV plots. The filters applied were as follows:

- Genotype is heterozygous in child (1/0) and homozygous in both parents (0/0)
- Child RD >20, Mother RD>20, Father RD>20
- Remove variants with >1 alternative read in either parent
- VAF>0.3 and VAF<0.7 for child
- Remove SNVs within 20 bp of each other. While this is likely removing true MNVs, the error mode was very high for clustered mutations.
- Removed DNMs if child RD >98 [13]
- Removed DNMs that fell within known segmental duplication regions as defined by
- UCSC (http://humanparalogy.gs.washington.edu/build37/data/ GRCh37GenomicSuperDup.tab)
- Removed DNMs that fell in highly repetitive regions (http://humanparalogy.gs.washington.edu/ build37/data/GRCh37simpleRepeat.txt)
- For DNM calls that fell on the X chromosome these slightly modified filters were used:
  - For DNMs that fell in PAR regions, the filters were unchanged from the autosomal calls apart from allowing for both heterozygous (1/0) and hemizygous (1) calls in males
  - For DNMs that fell in non-PAR regions the following filters were used:
    - For males: RD>20 in child, RD>20 in mother, no RD filter on father
    - For males: the genotype must be hemizygous (1) in child and homozygous in mother (0/0)
    - For females: RD>20 in child, RD>20 in mother, RD>10 in father

555
556

**DNM filtering and identifying hypermutated individuals in DDD**

To identify hypermutated individuals in the DDD study we started with exome sequencing data from the DDD study of families with a child with a severe, undiagnosed developmental disorder. The recruitment of these families has been described previously[62]: families were recruited at 24 clinical genetics centers within the UK National Health Service and the Republic of Ireland. Families gave informed consent to participate, and the study was

563  approved by the UK Research Ethics Committee (10/H0305/83, granted by the Cambridge
564  South Research Ethics Committee, and GEN/284/12, granted by the Republic of Ireland
565  Research Ethics Committee). Sequence alignment and variant calling of SNV and
566  insertions/deletions were conducted as previously described. De novo mutations were called
567  using DeNovoGear and filtered as previously[63,11]. The analysis in this paper was conducted
568  on a subset (7,930 parent offspring trios) of the full current cohort which was not available at
569  the start of this research.

571  In the DDD study, we identified 9 individuals out of 7,930 parent-offspring trios with an
572  increased number of exome DNMs after accounting for parental age (7-17 exome DNMs
573  compared to an expected number of ~2). These were subsequently submitted along with
574  their parents for PCR-free whole-genome sequencing at >30x mean coverage using
575  Illumina 150bp paired end reads and in house WSI sequencing pipelines. Reads were
576  mapped with bwa (v0.7.15)[64]. DNMs were called from these trios using DeNovoGear[63] and
577  were filtered as follows:
- Read depth (RD) of child > 10, mother RD > 10, father RD > 10
- Alternative allele read depth in child >2
- Filtered on strand bias across parents and child (p-value >0.001, Fisher's exact test)
- Removed DNMs that fell within known segmental duplication regions as defined by
- UCSC (http://humanparalogy.gs.washington.edu/build37/data/ GRCh37GenomicSuperDup.tab)
- Removed DNMs that fell in highly repetitive regions (http://humanparalogy.gs.washington.edu /build37/data/GRCh37simpleRepeat.txt)
- Allele frequency in gnomAD < 0.01
- VAF <0.1 for both parents
- Removed mutations if both parents have >1 read supporting the alternative allele
- Test to see if VAF in child is significantly greater than the error rate at that site as defined by error sites estimated using Shearwater [65].
- Posterior probability from DeNovoGear > 0.00781[63,11]
- Removed DNMs if child RD >200.

593  After applying these filters, this resulted in 1,367 DNMs. All of these DNMs were inspected in
594  the Integrative Genome Viewer[66] and removed if they appeared to be false positives. This
595  resulted in a final set of 916 DNMs across the 9 trios. One of the 9 had 277 dnSNVs genome
596  wide while the remaining had expected numbers (median number of 81 dnSNVs).

598  **Parental phasing of *de novo* mutations**
599  To phase the DNMs in both 100kGP and DDD we used a custom script which used the
600  following read-based approach to phase a DNM. This first searches for heterozygous
601  variants within 500 bp of the DNM that was able to be phased to a parent (so not
602  heterozygous in both parents and offspring). We then examined the reads or read pairs
603  which included both the variant and the DNM and counted how many times we observed the
604  DNM on the same haplotype of each parent. If the DNM appears exclusively on the same
605  haplotype as a single parent then that was determined to originate from that parent. We
606  discarded DNMs that had conflicting evidence from both parents. This code is available on
607  GitHub (https://github.com/queenjobo/PhaseMyDeNovo).

609  **Analysis of effect of parental age on germline mutation rate**
610  To assess the effect of parental age on germline mutation rate we ran the following
611  regressions on autosomal DNMs. On all (unphased) DNMs we ran two separate regressions

612 for SNVs and indels. We chose a negative Binomial generalized linear model here as the
613 Poisson was found to be overdispersed. We fitted the following model using a negative
614 Binomial GLM with an identity link where $Y$ is the number of DNMs for an individual:

$$E(Y) = \beta_0 + \beta_1 paternal\_age + \beta_2 maternal\_age$$

618 For the phased DNMs we fit the following two models using a negative Binomial GLM with
619 an identity link where $Y_{maternal}$ is the number of maternally derived DNMs and $Y_{paternal}$ is the
620 number of paternally derived DNMs:

$$E(Y_{paternal}) = \beta_0 + \beta_1 paternal\_age$$
$$E(Y_{maternal}) = \beta_0 + \beta_1 maternal\_age$$

**Identifying hypermutated individuals in 100kGP**

627 To identify hypermutated individuals in the 100kGP cohort we first wanted to regress out the
628 effect of parental age as described in the parental age analysis. We then looked at the
629 distribution of the studentized residuals and then, assuming these followed a $t$ distribution
630 with $N-3$ degrees of freedom, calculated a t-test p-value for each individual. We took the
631 same approach for the number of indels, except in this case $Y$ would be the number of de
632 novo indels.

634 We identified 21 individuals out of 12,471 parent-offspring trios with significantly increased
635 number of dnSNVs genome wide ($p < 0.05/12471$). We performed multiple quality control
636 analyses which included examining the mutations in the Integrative Genomics Browser for
637 these individuals to examine DNM calling accuracy, looking at the relative position of the
638 DNMs across the genome and examining the mutational spectra of the DNMs to identify any
639 well known sequencing error mutation types. We identified 12 that were not truly
640 hypermutated. The majority of false positives (10) were due to a parental somatic deletion in
641 blood increasing the number of apparent DNMs (Supplemental Figure 14). These individuals
642 had some of the highest number of DNMs called (up to 1379 DNMs per individual). For each
643 of these 10 individuals, the DNM calls all clustered to a specific region in a single
644 chromosome. In this same corresponding region in the parent, we observed a loss of
645 heterozygosity when calculating the heterozygous/homozygous ratio. In addition, many of
646 these calls appeared to be low level mosaic in that same parent. This type of event has
647 previously been shown to create artifacts in CNV calls and is referred to as a 'Loss of
648 Transmitted Allele' event[67]. The remaining 2 false positives were due to bad data quality in
649 either the offspring or one of the parents leading to poor DNM calls. The large number of
650 DNMs in these false positive individuals also led to significant underdispersion in the model
651 so after removing these 12 individuals we reran the regression model and subsequently
652 identified 11 individuals which appeared truly hypermutated ($p< 0.05/12,459$).

**Extraction of mutational signatures**

655 Mutational signatures were extracted from maternally and paternally phased autosomal
656 DNMs , 24 controls (randomly selected), 25 individuals (father with a cancer diagnosis prior
657 to conception), 27 individuals (mother with a cancer diagnosis prior to conception) and 12
658 hypermutated individuals that we identified. All DNMs were lifted over to GRCh37 prior to
659 signature extraction (100kGP samples are a mix of GRCh37 and GRCh38) and through the
660 liftover process a small number of 100kGP DNMs were lost (0.09% overall, 2 DNMs lost

661 across all hypermutated individuals). The mutation counts for all the samples can be found
662 in Supplemental Table 1. This was done using SigProfiler (v1.0.17) and these signatures are
663 extracted and subsequently mapped on to COSMIC mutational signatures (COSMIC v91,
664 Mutational Signature v3.1)[19,40]. Sigprofiler defaults to selecting a solution with higher
665 specificity than sensitivity. A solution with 4 de-novo signatures was chosen as optimal by
666 SigProfiler for the 12 hypermutator samples. Another stable solution with five de-novo
667 signatures was also manually deconvoluted, which has been considered as the final
668 solution. The mutation probability for mutational signature SBSHYP can be found in
669 Supplemental Table 2.

**Signature comparison to external exposures**

672 We compared the extracted signatures from these hypermutated individuals to a compilation
673 of previously identified signatures caused by enviromental mutagens from the literature. The
674 environmental signatures were compiled from Kucab et al (Cell 2019)[24], Pich et al (Nautre
675 Genetics 2019)[51] and Volkova et al (Nature Communications 2020)[52]. Comparison was
676 calculated as the cosine similarity between the different signatures.

**Defining set of genes involved in DNA repair**

679 We compiled a list of DNA repair genes which were taken from an updated version of the
680 table in Lange et al, Nature Reviews Cancer 2011
681 (https://www.mdanderson.org/documents/Labs/Wood- Laboratory/human-dna-repair-
682 genes.html)[68]. These can be found in Supplemental Table 3. These are annotated with the
683 pathways they are involved with (eg. nucleotide-excision repair, mismatch repair ). A 'rare'
684 variant is defined as those with an allele frequency of <0.001 for heterozygous variants and
685 those with an allele frequency of <0.01 for homozygous variants in both 1000 Genomes as
686 well as across the 100kGP cohort.

**Kinetic characterization of MPG**

689 The A135T variant of MPG was generated by site-directed mutagenesis and confirmed by
690 sequencing both strands. The catalytic domain of WT and A135T MPG were expressed in
691 BL21(DE3) Rosetta2 *E. coli* and purified as described for the full-length protein [69]. Protein
692 concentration was determined by absorbance at 280 nm. Active concentration was
693 determined by electrophoretic mobility shift assay with 5'FAM-labeled pyrolidine-DNA
694 (Supplemental Figure 8) [48]. Glycosylase assays were performed with 50 mM NaMOPS, pH
695 7.3, 172 mM potassium acetate, 1 mM DTT, 1 mM EDTA, 0.1 mg/mL BSA at 37 $^{\circ}$C. For
696 single turnover glycosylase activity, a 5'-FAM-labeled duplex was annealed by heating to 95
697 $^{\circ}$C and slowly cooling to 4 $^{\circ}$C (see Supplemental Figure 9). DNA substrate concentration
698 was varied between 10 and 50 nM and MPG concentration was maintained in at least 2-fold
699 excess over DNA from 25 to 10,000 nM. Timepoints were quenched in 0.2 M NaOH, heated
700 to 70 $^{\circ}$C for 12.5 min, then mixed with formamide/EDTA loading buffer and analyzed by 15%
701 denaturing polyacrylamide gel electrophoresis. Fluorescence was quantified with a Typhoon
702 5 imager and ImageQuant software (GE). The fraction of product was fit by a single
703 exponential equation to determine the observed single turnover rate constant ($k_{obs}$). For Hx
704 excision, the concentration dependence was fit by the equation $k_{obs} = k_{max} [E]/(K_{1/2}+[E])$, in
705 which the $K_{1/2}$ is the concentration at which half the maximal rate constant ($k_{max}$) was
706 obtained and [E] is the concentration of enzyme. It was not possible to measure the $K_{1/2}$ for
707 $\varepsilon$A excision using a fluorescence-based assay due to extremely tight binding [70]. Multiple
708 turnover glycosylase assays were performed with 5 nM MPG and 10–40-fold excess of
709 substrate (Supplemental Figure 9).

710

711 **Estimating the fraction of variance explained**

712 To estimate the fraction of germline mutation variance explained by several factors, we fit
713 the
714 following negative Binomial GLMs with an identity link. Data quality is likely to correlate with
715 the number of DNMs detected so to reduce this variation we used a subset of the 100kGP
716 dataset which had been filtered on some base quality control (QC) metrics by the
717 Bioinformatics team at GEL:

718 ● cross-contamination < 5%
719 ● mapping rate > 75%
720 ● mean sample coverage > 20
721 ● insert size <250

722 We then included the following variables to try and capture as much of the residual
723 measurement error which may also be impacting DNM calling. In brackets are the
724 corresponding variable names used in the models below:

725 ● Mean coverage for the child, mother and father (*child_mean_RD, mother_mean_RD,*
726 *father_mean_RD*)
727 ● Proportion of aligned reads for the child, mother and father (*child_prop_aligned,*
728 *mother_prop_aligned , father_prop_aligned*)
729 ● Number of SNVs called for child, mother and father (child_snvs, mother_snvs,
730 father_snvs)
731 ● Median VAF of DNMs called in child (*median_VAF*)
732 ● Median 'Bayes Factor' as outputted by Platypus for DNMs called in the child. This is
733 a metric of DNM quality (*median_BF*).

734

735 The first model only included parental age:

736

737 $$E(Y) = \beta_0 + \beta_1 paternal\_age + \beta_2 maternal\_age$$

738

739 The second model also included data quality variables as described above:

740

741 $$E(Y) = \beta_0 + \beta_1 paternal\_age + \beta_2 maternal\_age + \beta_3 child\_mean\_RD$$
742 $$+ \beta_4 \; mother\_mean\_RD + \beta_5 \; father\_mean\_RD + \beta_6 child\_prop\_aligned$$
743 $$+ \beta_7 mother\_prop\_aligned + \beta_8 \; father\_prop\_aligned + \beta_9 child\_snvs$$
744 $$+ \beta_{10} mother\_snvs + \beta_{11} \; father\_snvs + \beta_{12} median\_VAF + \beta_{13} median\_BF$$

745

746

747 The third model included a variable for excess mutations in the 11 confirmed hypermutated
748 individuals (hm_excess) in the 100kGP dataset. This variable was the total number of
749 mutations subtracted by the median number of DNMs in the cohort (65), $Y_{hypermutated} -$
750 median(Y) for these 11 individuals and 0 for all other individuals.

751

752 $$E(Y) = \beta_0 + \beta_1 paternal\_age + \beta_2 maternal\_age + \beta_3 child\_mean\_RD$$
753 $$+ \beta_4 mother\_mean\_RD + \beta_5 \; father\_mean\_RD + \beta_6 child\_prop\_aligned$$
754 $$+ \beta_7 mother\_prop\_aligned + \beta_8 \; father\_prop\_aligned + \beta_9 child\_snvs$$
755 $$+ \beta_{10} mother\_snvs + \beta_{11} \; father\_snvs + \beta_{12} median\_VAF + \beta_{13} median\_BF$$
756 $$+ \beta_{14} hm\_excess$$

757

758 The fraction of variance (*F*) explained after accounting for Poisson variance in the mutation
759 rate was calculated in a similar way to Kong et al using the following formula[3].

760
$$F = pseudoR^2 \frac{1 - \bar{Y}}{Var(Y)}$$

761

762 McFadden's pseudo $R^2$ was used here as a Negative binomial GLM was fitted. We repeated
763 these analyses fitting an ordinary least squares regression, as was done in Kong et al[3],
764 using the $R^2$ and got comparable results. To calculate a 95% confidence interval we used a
765 bootstrapping approach. We sampled with replacement 1,000 times and extracted the 2.5%
766 and 97.5% percentiles.

767

768 **Analysis of contribution of rare variants in DNA repair genes**
769 We fit 8 separate regressions to assess the contribution of rare variants in DNA repair genes
770 (compiled as described previously). These were across three different sets of genes:
771 variants in all DNA repair genes, variants in a subset of DNA repair genes known to be
772 associated with BER, MMR, NER or a DNA polymerase and variants within this subset that
773 have also been associated with a cancer phenotype. For this we downloaded all ClinVar
774 entries as of October 2019 and searched for germline 'pathogenic' or 'likely pathogenic'
775 variants annotated with cancer [54]. We tested both all nonsynonymous variants and just
776 protein truncating variants (PTVs) for each set. To assess the contribution of each of these
777 sets we created two binary variables per set indicating a presence or absence of a maternal
778 or paternal variant for each individual and then ran a negative binomial regression for each
779 subset including these as independent variables along with hypermutation status, parental
780 age and QC metrics as described in the previous section.

781

782 **Simulations to explore effect estimates of fraction of variance explained by paternal**
783 **age from downsampling**
784 To explore how the estimates of the fraction of variance of the number of DNMs is explained
785 by paternal age varies with downsampling we first simulated a random sample as follows
786 10,000 times:
- Randomly sample 78 trios (the number of trios in Kong et al.[3].)
- Fit OLS of $E(Y) = \beta_0 + \beta_1 paternal\_age$
- Estimated fraction of variance (*F*) as described in Kong et al[3].

790 We found that the median fraction explained was 0.77, sd of 0.13 and with 95% of
791 simulations fallings between 0.51 and 1.00.

792

793 **Identifying parents with cancer diagnosis prior to birth of offspring**
794 To identify parents who had received a cancer diagnosis prior to the conception of their child
795 we examined the admitted patient care hospital episode statistics of these parents. There
796 were no hospital episode statistics available prior to 1997 and many individuals did not have
797 any records until after the birth of the child. To ensure comparisons were not biased by this
798 we first subsetted to parents who had at least one episode statistic recorded at least two
799 years prior to the child's year of birth. Two years prior to the child's birth was our best
800 approximation for prior to conception without the exact child date of birth. This resulted in
801 2,891 fathers and 5,508 mothers. From this set we then extracted all entries with ICD10
802 codes with a "C" prefix which corresponds to malignant neoplasms and "Z85" which
803 corresponds to a personal history of malignant neoplasm. We defined a parent as having a
804 cancer diagnosis prior to conception if they had any of these codes recorded >=2 years prior

805  to the child's year of birth. We also extracted all entries with ICD10 code "Z511" which codes
806  for an 'encounter for antineoplastic chemotherapy and immunotherapy'.
807  Two fathers of hypermutated individuals who we suspect had chemotherapy prior to
808  conception did not meet these criteria as the father of GEL_5 received chemotherapy for
809  treatment for SLE and not cancer and for the father of GEL_8 the hospital record 'personal
810  history of malignant neoplasm' were entered after the conception of the child (Supplemental
811  Table 4).
812  To compare the number of dnSNVs between the group of individuals with parents with and
813  without cancer diagnoses we used a Wilcoxon test on the residuals from the negative
814  binomial regression on dnSNVs correcting for parental age, hypermutation status and data
815  quality. To look at the effect of maternal cancer on dnSNVs we matched these individuals on
816  maternal and paternal age with sampling replacement with 20 controls for each of the 27
817  individuals. We found a significant increase in DNMs (74 compared to 65 median dnSNVs, p
818  = 0.001, Wilcoxon Test).
819
820  **SNP heritability analysis**
821  For this analysis we started with the same subset of the 100kGP dataset that had been
822  filtered as described in the analysis on the impact of rare variants in DNA repair genes
823  across the cohort (see above). To ensure variant quality we subsetted to variants that have
824  been observed in genomes from gnomAD (v3) [71]. These were then filtered by ancestry to
825  parent-offspring trios where both the parents and child mapped on to the 1000 Genomes
826  GBR subpopulations. The first 10 principal components were subsequently included in the
827  heritability analyses. To remove cryptic relatedness we removed individuals with estimated
828  relatedness >0.025 (using GCTA grm-cutoff 0.025). This resulted in a set of 6,352 fathers
829  and 6,329 mothers. The phenotype in this analysis was defined as the residual from the
830  negative binomial regression of the number of DNMs after accounting for parental age,
831  hypermutation status and several data quality variables as described when estimating the
832  fraction of DNM count variation explained (see Methods above). To estimate heritability we
833  ran GCTA's GREML-LDMS on two LD stratifications and three MAF bins (0.001-0.01,0.01-
834  0.05,0.05-1).[55] For mothers this was run with the --reml-no-constrain option because
835  otherwise it would not converge. (Supplemental Table 8)
836

837

1. Conrad, D. F. *et al.* Variation in genome-wide mutation rates within and between human families. *Nat. Genet.* **43**, 712–714 (2011).

2. Roach, J. C. *et al.* Analysis of Genetic Inheritance in a Family Quartet by Whole-Genome Sequencing. *Science* vol. 328 636–639 (2010).

3. Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to disease risk. *Nature* **488**, 471–475 (2012).

4. Sasani, T. A. *et al.* Large, three-generation human families reveal post-zygotic mosaicism and variability in germline mutation accumulation. *Elife* **8**, (2019).

5. Lindsay, S. J., Rahbari, R., Kaplanis, J., Keane, T. & Hurles, M. E. Similarities and differences in patterns of germline mutation between mice and humans. *Nat. Commun.* **10**, 4053 (2019).

6. Harris, K. Evidence for recent, population-specific evolution of the human mutation rate. *Proceedings of the National Academy of Sciences* vol. 112 3439–3444 (2015).

7. Yang, S. *et al.* Parent–progeny sequencing indicates higher mutation rates in heterozygotes. *Nature* vol. 523 463–467 (2015).

8. Amos, W. Flanking heterozygosity influences the relative probability of different base substitutions in humans. *R Soc Open Sci* **6**, 191018 (2019).

9. Girard, S. L. *et al.* Paternal Age Explains a Major Portion of De Novo Germline Mutation Rate Variability in Healthy Individuals. *PLOS ONE* vol. 11 e0164212 (2016).

10. Mitra, I. *et al.* Genome-wide patterns of de novo tandem repeat mutations and their contribution to autism spectrum disorders. doi:10.1101/2020.03.04.974170.

11. Study, Deciphering Developmental Disorders. Prevalence and Architecture of De Novo Mutations in Developmental Disorders. *Nature* vol. 542 433–438 (2017).

12. Goldmann, J. M. *et al.* Parent-of-origin-specific signatures of de novo mutations. *Nat. Genet.* **48**, 935–939 (2016).

13. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat. Genet.* **48**, 126–133 (2016).

865  14. Gao, Z. *et al.* Overlooked roles of DNA damage and maternal age in generating human

866      germline mutations. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 9491–9500 (2019).

867  15. Jónsson, H. *et al.* Parental influence on human germline de novo mutations in 1,548

868      trios from Iceland. *Nature* **549**, 519–522 (2017).

869  16. Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R.

870      Deciphering Signatures of Mutational Processes Operative in Human Cancer. *Cell*

871      *Reports* vol. 3 246–259 (2013).

872  17. Nik-Zainal, S. *et al.* 604 Cancer Genomics, Epigenetics and Genomic Instability.

873      Mutational Processes Shaping the Genomes of Twenty-one Breast Cancers. *European*

874      *Journal of Cancer* vol. 48 S144 (2012).

875  18. Phillips, D. H. Mutational spectra and mutational signatures: Insights into cancer

876      aetiology and mechanisms of DNA damage and repair. *DNA Repair* **71**, 6–11 (2018).

877  19. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature*

878      **578**, 94–101 (2020).

879  20. Demanelis, K. *et al.* Determinants of telomere length across human tissues. *Science*

880      **369**, (2020).

881  21. Moore, L., Cagan, A., Coorens, T., Neville, M. D. C. & Sanghvi, R. The mutational

882      landscape of human somatic and germline cells. *bioRxiv* (2020).

883  22. Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat.*

884      *Genet.* **47**, 1402–1407 (2015).

885  23. Stratton, M. R., Campbell, P. J. & Andrew Futreal, P. The cancer genome. *Nature* vol.

886      458 719–724 (2009).

887  24. Kucab, J. E. *et al.* A Compendium of Mutational Signatures of Environmental Agents.

888      *Cell* **177**, 821–836.e16 (2019).

889  25. Tawn, E. J. & Janet Tawn, E. Hereditary Effects of Radiation: UNSCEAR 2001 Report

890      to the General Assembly, with Scientific Annex. *Journal of Radiological Protection* vol.

891      22 121–122 (2002).

892  26. Forster, L., Forster, P., Lutz-Bonengel, S., Willkomm, H. & Brinkmann, B. Natural

893    radioactivity and human mitochondrial DNA mutations. *Proc. Natl. Acad. Sci. U. S. A.*

894    **99**, 13950–13954 (2002).

895    27. Adewoye, A. B., Lindsay, S. J., Dubrova, Y. E. & Hurles, M. E. The genome-wide effects

896        of ionizing radiation on mutation induction in the mammalian germline. *Nat. Commun.* **6**,

897        6684 (2015).

898    28. Huang, K.-L. *et al.* Pathogenic Germline Variants in 10,389 Adult Cancers. *Cell* **173**,

899        355–370.e14 (2018).

900    29. Dutil, J. *et al.* Germline variants in cancer genes in high-risk non-BRCA patients from

901        Puerto Rico. *Sci. Rep.* **9**, 17769 (2019).

902    30. Hu, C. *et al.* Association Between Inherited Germline Mutations in Cancer

903        Predisposition Genes and Risk of Pancreatic Cancer. *JAMA* vol. 319 2401 (2018).

904    31. Huang, Q.-Y. *et al.* Mutation patterns at dinucleotide microsatellite loci in humans. *Am.*

905        *J. Hum. Genet.* **70**, 625–634 (2002).

906    32. Gymrek, M., Willems, T., Reich, D. & Erlich, Y. Interpreting short tandem repeat

907        variations in humans using mutational constraint. *Nat. Genet.* **49**, 1495–1501 (2017).

908    33. Sun, J. X. *et al.* A direct characterization of human mutation based on microsatellites.

909        *Nature Genetics* vol. 44 1161–1165 (2012).

910    34. Monckton, D. G. *et al.* Minisatellite mutation rate variation associated with a flanking

911        DNA sequence polymorphism. *Nat. Genet.* **8**, 162–170 (1994).

912    35. Crowley, J. J. *et al.* Common-variant associations with fragile X syndrome. *Mol.*

913        *Psychiatry* **24**, 338–344 (2019).

914    36. Kaplanis, J. *et al.* Evidence for 28 genetic disorders discovered by combining healthcare

915        and research data. *Nature* **586**, 757–762 (2020).

916    37. Uchimura, A. *et al.* Germline mutation rates and the long-term phenotypic effects of

917        mutation accumulation in wild-type laboratory mice and mutator mice. *Genome Res.* **25**,

918        1125–1134 (2015).

919    38. Cawthon, R. M. *et al.* Germline mutation rates in young adults predict longevity and

920        reproductive lifespan. *Sci. Rep.* **10**, 10001 (2020).

921    39. Liu, P. *et al.* An Organismal CNV Mutator Phenotype Restricted to Early Human

922        Development. *Cell* **168**, 830–842.e7 (2017).

923    40. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic*

924        *Acids Res.* **47**, D941–D947 (2019).

925    41. Lehmann, A. R., McGibbon, D. & Stefanini, M. Xeroderma pigmentosum. *Orphanet J.*

926        *Rare Dis.* **6**, 70 (2011).

927    42. Pippard, E. C., Hall, A. J., Barker, D. J. & Bridges, B. A. Cancer in homozygotes and

928        heterozygotes of ataxia-telangiectasia and xeroderma pigmentosum in Britain. *Cancer*

929        *Res.* **48**, 2929–2932 (1988).

930    43. Jager, M. *et al.* Deficiency of nucleotide excision repair is associated with mutational

931        signature observed in cancer. *Genome Res.* **29**, 1067–1077 (2019).

932    44. Miccoli, L. *et al.* The combined effects of xeroderma pigmentosum C deficiency and

933        mutagens on mutation rates in the mouse germ line. *Cancer Res.* **67**, 4695–4699

934        (2007).

935    45. Laskowski, R. A., Stephenson, J. D., Sillitoe, I., Orengo, C. A. & Thornton, J. M. VarSite:

936        Disease variants and protein structure. *Protein Sci.* **29**, 111–119 (2020).

937    46. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting

938        the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **47**,

939        D886–D894 (2019).

940    47. Glassner, B. J., Rasmussen, L. J., Najarian, M. T., Posnick, L. M. & Samson, L. D.

941        Generation of a strong mutator phenotype in yeast by imbalanced base excision repair.

942        *Proc. Natl. Acad. Sci. U. S. A.* **95**, 9997–10002 (1998).

943    48. Eyler, D. E., Burnham, K. A., Wilson, T. E. & O'Brien, P. J. Mechanisms of glycosylase

944        induced genomic instability. *PLoS One* **12**, e0174041 (2017).

945    49. Connor, E. E., Wilson, J. J. & Wyatt, M. D. Effects of substrate specificity on initiating

946        the base excision repair of N-methylpurines by variant human 3-methyladenine DNA

947        glycosylases. *Chem. Res. Toxicol.* **18**, 87–94 (2005).

948    50. Boot, A. *et al.* In-depth characterization of the cisplatin mutational signature in human

949      cell lines and in esophageal and liver tumors. *Genome Res.* **28**, 654–665 (2018).

950    51. Pich, O. *et al.* The mutational footprints of cancer therapies. *Nat. Genet.* **51**, 1732–1740

951      (2019).

952    52. Volkova, N. V. *et al.* Mutational signatures are jointly shaped by DNA damage and

953      repair. *Nat. Commun.* **11**, 2169 (2020).

954    53. Hock, H. & Shimamura, A. ETV6 in hematopoiesis and leukemia predisposition. *Semin.*

955      *Hematol.* **54**, 98–104 (2017).

956    54. Landrum, M. J. *et al.* ClinVar: public archive of interpretations of clinically relevant

957      variants. *Nucleic Acids Res.* **44**, D862–8 (2016).

958    55. Yang, J. *et al.* Genetic variance estimation with imputed variants finds negligible missing

959      heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).

960    56. Kilpivaara, O. & Aaltonen, L. A. Diagnostic cancer genome sequencing and the

961      contribution of germline variants. *Science* **339**, 1559–1562 (2013).

962    57. Wong, E. *et al.* Mbd4 inactivation increases C->T transition mutations and promotes

963      gastrointestinal tumor formation. *Proceedings of the National Academy of Sciences* vol.

964      99 14937–14942 (2002).

965    58. Bart, J. *et al.* An oncological view on the blood-testis barrier. *Lancet Oncol.* **3**, 357–363

966      (2002).

967    59. Wheway, G., Mitchison, H. M. & Genomics England Research Consortium.

968      Opportunities and Challenges for Molecular Understanding of Ciliopathies-The 100,000

969      Genomes Project. *Front. Genet.* **10**, 127 (2019).

970    60. Ouwehand, W. H., on behalf of the NIHR BioResource and the & 000 Genomes Project.

971      Whole-genome sequencing of rare disease patients in a national healthcare system.

972      doi:10.1101/507244.

973    61. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for

974      calling variants in clinical sequencing applications. *Nat. Genet.* **46**, 912–918 (2014).

975    62. Wright, C. F. *et al.* Genetic diagnosis of developmental disorders in the DDD study: a

976      scalable analysis of genome-wide research data. *Lancet* **385**, 1305–1314 (2015).

977    63.  Ramu, A. *et al.* DeNovoGear: de novo indel and point mutation discovery and phasing.

978         *Nat. Methods* **10**, 985–987 (2013).

979    64.  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler

980         transform. *Bioinformatics* **25**, 1754–1760 (2009).

981    65.  Gerstung, M., Papaemmanuil, E. & Campbell, P. J. Subclonal variant calling with

982         multiple samples and prior knowledge. *Bioinformatics* **30**, 1198–1204 (2014).

983    66.  Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV):

984         high-performance genomics data visualization and exploration. *Briefings in*

985         *Bioinformatics* vol. 14 178–192 (2013).

986    67.  Redon, R. *et al.* Global variation in copy number in the human genome. *Nature* **444**,

987         444–454 (2006).

988    68.  Lange, S. S., Takata, K.-I. & Wood, R. D. DNA polymerases and cancer. *Nat. Rev.*

989         *Cancer* **11**, 96–110 (2011).

990    69.  Zhang, Y. & O'Brien, P. J. Repair of Alkylation Damage in Eukaryotic Chromatin

991         Depends on Searching Ability of Alkyladenine DNA Glycosylase. *ACS Chem. Biol.* **10**,

992         2606–2615 (2015).

993    70.  Wolfe, A. E. & O'Brien, P. J. Kinetic mechanism for the flipping and excision of 1,N(6)-

994         ethenoadenine by human alkyladenine DNA glycosylase. *Biochemistry* **48**, 11357–

995         11369 (2009).

996    71.  Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in

997         141,456 humans. *Nature* **581**, 434–443 (2020).

998