

1 **Clonal myelopoiesis in the UK Biobank cohort: *ASXL1* mutations are strongly**
2 **associated with smoking**

3
4
5 Ahmed A.Z. Dawoud¹, William J Tapper¹ & Nicholas C.P. Cross^{1,2}

6
7
8
9 ¹ Faculty of Medicine, University of Southampton, Southampton, UK

10 ² Wessex Regional Genetics Laboratory, Salisbury NHS Foundation Trust, Salisbury, UK

11
12
13
14
15
16
17 Correspondence to:

18
19 Professor N.C.P. Cross
20 Wessex Regional Genetics Laboratory
21 Salisbury NHS Foundation Trust
22 Salisbury SP2 8BJ, UK

23
24 Tel: +(44) 1722 429080

25 Fax: +(44) 1722 331531

26 email: ncpc@soton.ac.uk

27

28

29

30 **Abstract**

31

32 We sought to determine the significance of myeloid clonal hematopoiesis (CH) in the UK
33 Biobank cohort (n=502,524, median age=58 years). Utilizing SNP array (n=486,941) and
34 whole exome sequencing data (n=49,956), we identified 1166 participants with myeloid CH,
35 defined by myeloid-associated mosaic chromosome abnormalities (mCA) and/or likely
36 somatic driver mutations in *DNMT3A*, *TET2*, *ASXL1*, *JAK2*, *SRSF2* or *PPM1D*. Myeloid CH
37 increased by 1.1-fold per annum (myeloid mCA, $P=1.57 \times 10^{-38}$; driver mutations, $P=5.89 \times 10^{-47}$).
38 Genome-wide association analysis identified two distinct signals within *TERT* that
39 predisposed to myeloid CH, plus a weaker signal corresponding to the *JAK2* 46/1 haplotype.
40 Specific subtypes of myeloid CH were associated with several blood features and clinical
41 phenotypes, including *TET2* mutations and chronic obstructive pulmonary disease. Smoking
42 history was significantly associated with myeloid CH: 53% of myeloid CH cases were smokers
43 compared to 44% of controls ($P=3.38 \times 10^{-6}$), a difference principally due to current (OR=1.10;
44 $P=6.14 \times 10^{-6}$) rather than past smoking ($P=0.08$). Breakdown of CH by specific mutation type
45 revealed that *ASXL1* loss of function mutations were most strongly associated with current
46 smoking status (OR=1.07; $P=1.92 \times 10^{-5}$), and the only abnormality associated with past
47 smoking (OR=1.04; $P=0.0026$). We suggest that the inflammatory environment induced by
48 smoking may promote the outgrowth of *ASXL1*-mutant clones.

49

50

51

52 **Running title:** *ASXL1* mutations are associated with smoking history

53 **Keywords:** *ASXL1*, smoking, UK Biobank, clonal hematopoiesis

54

55 **Introduction**

56

57 Age-related clonal hematopoiesis (CH) is a widespread phenomenon that predisposes to the
58 development of hematological malignancies as well as some non-hematological conditions.
59 First suggested more than 20 years ago by skewed X-chromosome inactivation patterns in
60 myeloid cells of elderly females,¹ clonality in the absence of a hematological neoplasm was
61 proven by the finding of acquired *TET2* driver mutations in a subset of skewed cases.² CH as
62 a more general phenomenon was defined by several lines of evidence, including (i) the
63 prevalence of *JAK2* V617F was much greater than myeloproliferative neoplasms (MPN) in
64 randomly selected subjects undergoing hospital-based clinical investigations;³ (ii) mosaic
65 chromosome alterations (mCA), indicative of myeloid or lymphoid clonality, were seen in
66 large cohorts ascertained for non-hematological conditions^{4, 5} and (iii) myeloid malignancy-
67 associated driver mutations, most commonly involving *DNMT3A*, *TET2* or *ASXL1*, were found
68 at a much higher frequency than expected in similar large, non-hematological cohorts.⁶⁻⁸

69

70 The studies above determined that CH is strongly age dependent, with >10% of individuals
71 aged >65-70 years of age being affected, but only 1% of those <50 years old. Subsequent
72 analysis using more sensitive mutation detection methods indicated that driver mutation-
73 associated clonal expansion is more common⁹ than previously thought and indeed nearly
74 ubiquitous when the limit of detection is pushed to as low as 0.03% variant allele frequency
75 (VAF).¹⁰ The picture is further complicated by the fact that whole exome sequencing (WES)
76 and whole genome sequencing (WGS) studies have found that CH often predominates in the
77 absence of discernible driver mutations, potentially as a consequence of stochastic drift
78 acting on a small population of active hematopoietic stem cells.⁶⁻⁸

79

80 Although individuals with CH have a 10-fold excess risk of developing a hematological
81 neoplasm, the actual rate of progression is only 1-2% per annum, and thus the term CH of
82 indeterminate potential (CHIP) is widely used.¹¹ However CH is also associated with all-cause
83 mortality, and specifically an increased risk of atherosclerotic cardiovascular disease after

84 adjusting for known risk factors.¹² CH has also been associated with other disorders, for
85 example chronic obstructive pulmonary disorders.¹³

86

87 UK Biobank is an ongoing prospective cohort study that aims to provide new insights into
88 the causes of common and uncommon diseases. Recruitment of approximately 500,000
89 community-dwelling, generally healthy subjects aged 40-69 years took place between 2006
90 and 2010 across the UK. Participants provided comprehensive demographic, psychosocial
91 and

92 medical information during an initial visit to a UK Biobank assessment centre along with
93 baseline blood and urine samples for genomic, biochemical and other laboratory
94 assessments. Long term follow-up is provided via linked medical records.¹⁴ In this study we
95 aimed to assess the prevalence, impact and causes of myeloid CH, defined by the presence
96 of myeloid mCA and/or myeloid malignancy-associated mutations, in the UK Biobank
97 cohort.

98 **Methods**

99 *Cohort structure.* UK Biobank participants were split into four phenotypic groups: myeloid
100 disorders, lymphoid disorders, other cancers, and cancer free based on the International
101 Classification of Disease codes (ICD version 10) recorded by the national cancer registry
102 (Data-Field 40006) and reason for admission to hospital (Data-Fields 41202, 41204, and
103 41270; Supplementary Table 1). In subsequent analyses, these four groups are collectively
104 referred to as the phenotypic groups. The ICD-10 coding system was introduced into the
105 national cancer registry in 1995 and to the hospital admissions in 1996. Consequently, any
106 relevant events that occurred before 1995/6 were not captured. The cancer registry data
107 was accessed as of July 2018 (most recent record December 2016); other clinical and
108 phenotype data was accessed as of August 2019 (most recent record February 2018), thus
109 providing data for a median of 9.1 years after recruitment and blood sampling (median age
110 at recruitment = 58 years). The four phenotypic groups were defined by events that
111 occurred at any time during the study period (1995 to 2018). All participants provided
112 informed consent according to the Declaration of Helsinki, and UK Biobank received ethical
113 approval from the North West multi-centre Research Ethics Committee (REC reference
114 11/NW/0382). The current study was conducted under approved UK Biobank application
115 number 35273.

116

117 *SNP genotype data.* Genotypic data for 488,377 samples were obtained from UK Biobank.
118 These samples were genotyped at Affymetrix laboratories (Santa Clara, USA) in batches of
119 approximately 4,700 samples using two similar microarrays, the UK BiLEVE array (807,411
120 markers) (n=49,950 samples) and the remainder by UK Biobank axiom array (825,927
121 markers). The two microarrays share 95% of SNPs and are thus highly comparable.
122 Routine quality control (QC) was performed by Affymetrix, which excluded poorly
123 genotyped SNPs, and multi-allelic SNPs (n= 35,014 across the two platform). Further QC was
124 performed by UK Biobank (Supplementary Methods) that restricted the released data to
125 805,426 markers, and 486,941 samples with no quality flags.¹⁵ Of the 488,377 participants
126 with available SNP data, we excluded a further 1,436 cases due to one or more of poor
127 genotyping quality (missingness above 5%; n=229), outlying levels of heterozygosity that

128 could not be explained by admixture or consanguinity (principal component-adjusted
129 heterozygosity above the mean 0.1903, n=744), gender mismatch (n=373), withdrawal of
130 consent (n=85), or absence of phenotypic data (n=10).

131

132 *Calling mosaic chromosomal alterations from SNP array data.* Raw input files were
133 generated for each sample which contained the B-allele frequency (BAF: the ratio of
134 intensity values for the A and B allele for each SNP in a single sample) and \log_2 R ratio (LRR:
135 the logarithm value to the base 2 of the ratio of the observed intensity to the expected
136 intensity for the diploid genome) for each SNP that passed QC. Regions of allelic imbalance
137 (AI) were detected using BAF segmentation to analyse the raw input files using the default
138 parameters for Affymetrix array data.¹⁶ To select likely somatic events and exclude potential
139 false positives, a custom script (see Supplementary Methods) was used to process the
140 resulting AI regions as follows. First, bedtools was used to merge AI regions with a minimum
141 density of 1 SNP per 20Kb that were separated by less than 2Mb. AI regions were then
142 scored based on the product of SNP density, heterozygosity rate and coverage (see
143 Supplementary Table 2). Finally, AI regions greater than 2Mb in total length with confidence
144 scores above an empirically defined threshold (≥ 9) were defined as mCA. These were
145 further broken down into copy number loss (CNL), copy number gain (CNG) or acquired
146 uniparental disomy (aUPD) based on the regions median LRR value. The parameters used
147 have been estimated to identify clonal fractions larger than 0.1, 0.2 and 0.27 for aUPD, CNL
148 and CNG, respectively.¹⁶ Full details are given in the Supplementary Methods. Since mCA
149 may be derived from myeloid or lymphoid cells, we correlated mCA with clinical phenotype
150 to specifically define myeloid mCA (see Results).

151

152 *Whole exome sequencing data.* Whole exome sequencing (WES) was performed by
153 Regeneron Genetics. In brief, read sequences were aligned to the reference human genome
154 (version GRCH38) using BWA-MEM, duplicate reads were marked using Picard tools, and
155 patient level genomic variant call files (gVCF) were generated using the WeCall variant caller
156 with ≥ 2 alternative reads to call a variant as described.¹⁷ Of the 49,996 subjects that were

157 successfully sequenced, 40 samples flagged by Regeneron Genetics for QC issues were
158 excluded from the analysis.¹⁷

159

160 *Identification of candidate somatic driver variants from WES data.* The gVCF files from UK
161 Biobank were converted to VCF format and filtered to remove variants with low read depth
162 (DP; <7 for SNVs, <10 for indels). SAMtools/Bcftools¹⁸ was used to merge the separate files
163 into one multi-sample VCF, split multi-allelic positions into separate variants and to
164 normalise the location of indels using their left most position.¹⁸ The combined VCF file was
165 annotated using Annovar and the RefSeq gene database¹⁹. Putative somatic mutations
166 (regardless of VAF) were identified in six genes known to be associated with myeloid
167 neoplasia that were exonic, had an alternate allele frequency $\leq 1\%$ in public databases of
168 common variation (1000 Genome, ExAC, ESP6500, gnomAD) and were either loss of function
169 (LOF) mutations (*TET2*, *DNMT3A*, *ASXL1*, *PPM1D*) or known oncogenic hotspots (*DNMT3A*
170 R882, *JAK2* V617F, *SRSF2* P95; see Supplementary Methods).⁷

171

172 *Statistical analyses:*

173 *Association of mCA categories and somatic driver mutations with phenotypic group.* The
174 frequency of mCA events, each of its subcategories (aUPD, CNG or CNL) and somatic driver
175 mutations were tested for association with either myeloid, lymphoid, or other cancers
176 compared to cancer free controls using Fisher's exact tests in SPSS (Version 25). The average
177 number of mCA events per sample in either myeloid, lymphoid, or other cancers were
178 compared against cancer free controls using Mann-Whitney U tests²⁰.

179

180 *Association of specific mCAs with hematological phenotype.* Autosomal mCAs were stratified
181 by type (aUPD, CNL, CNG), chromosome arm (*p* or *q*), and position (telomeric or interstitial).
182 Each type with at least one observation was tested for association with a hematological
183 phenotype (either myeloid, or lymphoid) in comparison to cancer free controls using
184 Fisher's exact tests in SPSS (version 25). A total of 416 specific mCAs were tested after
185 selecting mCA types with at least one observation. Interstitial events that were associated
186 with a myeloid phenotype and not previously recognised as a recurrent abnormality were

187 manually reviewed (see Supplementary Methods). After review, events associated with a
188 myeloid phenotype were grouped and hereafter referred to as myeloid mCAs.

189

190 *Regression of mCA against age.* The relationship between mCA and age was tested using
191 multivariable logistic regression in SPSS where mCA status was treated as the dependent,
192 age as a predictor and including gender as an independent covariate. This analysis was
193 repeated for each subcategory of mCA (aUPD, CNG or CNL), myeloid mCAs and myeloid
194 somatic mutations. The effect sizes were reported as odds ratios (OR) with 95% confidence
195 intervals (CI).

196

197 *The association of common variants with clonality.* Samples with SNP array data were split
198 into cases and controls. Cases were defined by the presence of one or more feature
199 associated with myeloid CH; either a myeloid mCA or at least one putative somatic mutation
200 in six driver genes (*JAK2 V617F*, *SRSF2 P95*, *DNMT3A R882* or any frameshift/stopgain
201 mutation in *DNMT3A*, *TET2*, *ASXL1* or *PPM1D*). Controls were defined as samples without
202 mCA, without likely somatic mutations in the genes of interest (including nonsynonymous
203 variants) and without evidence of any hematological malignancy during the study period. A
204 total of 265,112 SNPs with minor allele frequencies (MAF) greater than 10% and without
205 deviation from Hardy-Weinberg equilibrium ($P > 0.001$) were assessed for association with
206 myeloid CH using allelic chi-square tests to compare allele frequencies between cases and
207 controls. Association tests were performed using Plink V1.9.²¹ These results were visualised
208 using the qqman, qqnorm and qqplot plot procedures in R to generate a Manhattan plot
209 and quantile-quantile plot. In regions with multiple SNPs reaching genome-wide
210 significance, stepwise logistic regression was used to determine the number of independent
211 signals. All SNPs with $P < 10^{-8}$ and within 500kb of the index SNP were added to the
212 regression model in order of significance.

213

214 *The association of clonal hematopoiesis with smoking (past, current, any), clinical phenotype*
215 *blood traits and biochemistry.* We used the PHENome Scan ANalysis Tool (PHESANT)²² to test
216 the effect of myeloid CH, which is defined by myeloid mCAs and/or somatic mutations

217 associated with myeloid diseases, on selected phenotypes from the UK Biobank.
218 Phenotypes were tested using either ordinal (current and previous smoking), multinomial
219 (combined smoking status), logistic (clinical phenotypes, n=395), or linear regression (blood
220 features, n=29 or biochemical markers, n=30). All regressions included covariates for age
221 and sex with the addition of smoking for the analysis of clinical phenotypes and blood
222 features. Where appropriate, inverse normal transformation was applied to counteract
223 departures from normality. Clinical phenotypes were determined using the primary/main
224 diagnoses from hospital inpatient records (Data-field 41202) and limited to phenotypes
225 identified in $\geq 0.1\%$ of the UK Biobank participants.

226

227 *Survival analyses:* To test the association between myeloid CH and either all-cause
228 mortality, myocardial infarction (MI) or stroke, we used the survival package²³ in R to
229 perform Cox regression analyses with correction for age, sex and smoking status. Follow-up
230 times were calculated using the “lubridate”²⁴ package to determine the duration between
231 study entry to last registration in either the date of death (Data-field 4000) date of MI (Data-
232 field 42000) or date of stroke (Data-field 42006). Participants that had an event before the
233 date of entry were excluded.

234

235 Where relevant, P values were corrected for multiple testing using the false discovery rate
236 (FDR) and denoted P_{FDR} . Detailed methods plus associated figures and tables are provided in
237 the Supplementary Methods.

238 **Results**

239

240 *The study cohort.* The phenotypic breakdown of the UK Biobank cohort is summarised in
241 Table 1, along with cases for whom SNP array and WES data were available. The
242 classification shows the expected excess of myeloid malignancies in males (1.3 : 1).²⁵

243

244 *Identification of mosaic chromosomal alterations.* We analysed genome wide SNP array data
245 to identify autosomal regions of AI in all participants for whom the array data passed QC
246 (n=486,941), and X imbalances for female participants (n=264,083). In the absence of
247 matched constitutional DNA to definitively identify somatic events, blood cell clonality was
248 inferred by using an upper mBAF threshold to exclude events that were likely to be
249 constitutional. In total, our method identified 8,203 mCA >2Mb in size in 5,040 participants
250 (1% of 486,941 analysed samples; Supplementary Table 2) which broke down into aUPD
251 (n=4224), CNG (n=659) or CNL (n=3320). The frequency of these events (mCA, aUPD, CNG
252 and CNL) were all significantly higher in the myeloid and lymphoid disease groups versus
253 cancer free controls while only CNL were more frequent in other types of cancer versus
254 controls (Table 2). The incidence of mCAs was highest in the myeloid disorders group with
255 11% of samples affected (210/1913). Of these, more than 75% involved aUPD (158/210)
256 which corresponds to a 16-fold enrichment versus cancer free controls and the most
257 significantly associated mCA subcategory with disease (OR=16.39; $P_{FDR}=8.78 \times 10^{-124}$; Fisher's
258 exact test). The frequency of mCA in lymphoid disorders was much lower than the myeloid
259 group (363/12546; 2.9%) but the average number of events per positive sample (2.1) was
260 significantly higher compared to cancer free controls (1.6, $P=4.1 \times 10^{-6}$) or myeloid samples
261 (1.5, $P=0.015$) according to the Mann-Whitney U tests (Table 2).

262

263 The frequency of mCA increased with age and ranged between 0.85% at 40-45 years to
264 1.29% at 66-70 years over all UK Biobank participants with array data passing QC
265 (n=486,941) (Figure 1a). Using logistic regression, this corresponded to an annual increased
266 risk of acquiring a mCA by 1.02 fold after adjusting for gender (OR=1.02, $P=1.80 \times 10^{-19}$). The

267 risk of acquiring each subcategory of mCA also increased with age; aUPD (OR = 1.02,
268 $P=3.14 \times 10^{-14}$), CNG (OR = 1.04, $P=1.09 \times 10^{-6}$) and CNL (OR = 1.01, $P=3.66 \times 10^{-4}$).

269

270 *Association of specific mCA with hematological phenotype.* As expected, distinct mCAs were
271 associated with myeloid and lymphoid disorders. After excluding interstitial events that
272 failed manual review, a total of 17 mCAs involving 15 chromosomal arms were found to be
273 associated with myeloid disorders using Fisher's exact tests (Table 3). Hereafter, these 17
274 abnormalities are collectively referred to as myeloid mCAs. Of these, 9p aUPD was the most
275 significant individual abnormality associated with myeloid disorders (OR=2858,
276 $P_{FDR}=6.28 \times 10^{-191}$, Table 3). The risk of acquiring at least one of the myeloid mCAs was more
277 strongly associated with age (OR=1.10, $P=1.57 \times 10^{-38}$) than mCAs involving all chromosomes
278 (Figure 1b). Interstitial CNL of chromosome 13 was the most significant abnormality
279 associated with lymphoid disorders (OR=23.16; $P_{FDR}=1.24 \times 10^{-63}$, Supplementary Table 3).

280

281 *Clonality defined by somatic mutations in individual genes.* We restricted our analysis of the
282 WES data to a set of genes/mutations that collectively account for approximately 95% of
283 myeloid CH events in 6 genes (*DNMT3A*, *TET2*, *ASXL1*, *JAK2*, *SRSF2*, and *PPM1D*)^{6-8, 26}. As
284 summarised in Table 4, we identified 721 candidate driver mutations in 678 subjects, with
285 *DNMT3A* being the most commonly affected gene. Only 37 cases had more than one
286 variant, with frequency of myeloid disorders being higher (11.8%) in participants with more
287 than one variant compared to those with a single variant (5.8%). Of the 678 participants
288 with CH defined by somatic mutations, 18 (*JAK2*, n=11; *DNMT3A*, n=4; *ASXL1*, *PPM1D*, *TET2*,
289 n=1 of each) also had CH defined by mCA. A detailed list of variants in the 6 genes are given
290 in Supplemental Table 4.

291

292 As expected, the prevalence of these putative somatic mutations was shown to be greatest
293 in cases with myeloid disorders (22.5% versus 1.2% for cancer-free controls, $P_{FDR}=5.83 \times 10^{-38}$,
294 OR=23.7) compared to other groups (2.1% for lymphoid disorders versus cancer-free
295 controls; $P_{FDR}=0.02$; OR=1.7) using Fisher's exact tests. Looking at individual genes, the only

296 exception was *JAK2 V617F* which was most commonly seen in myeloid disorders. There was
297 less difference in the prevalence of myeloid mutations in participants who had non-
298 hematological cancers versus cancer-free controls (1.4% vs 1.2%; $P_{FDR}=0.04$; OR=1.2).

299

300 The frequency of CH defined by putative somatic mutations also increased with age and
301 ranged between 0.4% at age <45 years to 2.8% at age >65 years (Figure 1e). The risk of
302 acquiring any one of these mutations increased by 1.1 fold per year, which is the same as
303 the effect determined for myeloid mCA (OR=1.1, $P=5.89 \times 10^{-47}$). The six genes were also
304 significant when tested separately (Supplementary Table 5). The risk of acquiring a somatic
305 mutation in *SRSF2* had the largest annual increase with age at 1.2 fold ($P=3.66 \times 10^{-4}$)
306 although these mutations were the rarest (n=20) and absent in participants younger than 50
307 years old (Figure 1e).

308

309 *The association of CH with constitutional genetic variation.* A previous study has associated
310 germline variation at the *TERT* locus (rs34002450) with CH in the Icelandic population.⁷ To
311 examine the influence of genetic variation on myeloid CH in the UK Biobank cohort, we
312 performed a genome wide association study (GWAS) to assess the influence of common
313 polymorphisms in the 1,166 participants with myeloid CH defined by the presence of (i) any
314 myeloid malignancy-associated mCA (Table 3) and/or (ii) somatic mutations in the 6 genes
315 of interest. A total of 265,112 SNPs were tested for association with myeloid CH using allelic
316 chi-square tests to compare these cases with 30,892 controls that were free of any clonal
317 marker, or any hematological malignancy (Figure 2a,b). Three SNPs with genome-wide
318 significance were identified in the *TERT* gene. Two of these were associated with an
319 increased risk of developing CH (rs2853677 intron 2, OR=1.32, $P=5.6 \times 10^{-11}$; rs7726159
320 intron 3, OR=1.33, $P=4.2 \times 10^{-11}$) while the third and most significant single SNP was
321 protective (rs2736100 intron 2, OR=0.74, $P=3.1 \times 10^{-12}$) (Figure 2c; Supplementary Table 6).
322 Two of these SNPs (rs7726159, A allele, OR=1.19, $P=0.003$ and rs2853677, G allele,
323 OR=1.18, $P=0.004$) were identified as independent association signals using stepwise logistic
324 regression with an additive model and treating all three SNPs as covariates. A second signal

325 was seen just below the level of genome wide significance (Figure 2b) and included
326 rs3780381, rs17425819 and rs10974944. These SNPs are within *JAK2* and are in linkage
327 disequilibrium (LD) with the 46/1 haplotype, previously shown to be strongly associated
328 with acquisition of *JAK2* V617F.²⁷ This association signal disappeared when cases with *JAK2*
329 V617F (n=40) and mCA including *JAK2* (n=115) were removed from the analysis.

330

331 *The association between clonal hematopoiesis and smoking.* To examine the relationship
332 between smoking and CH, we used PHESANT to perform regression analyses of past, current
333 and combined smoking status in 32,058 participants consisting of 1,166 with myeloid CH
334 and 30,892 without myeloid CH and without any hematological malignancy. The odds of
335 having ever smoked (combined status) were significantly higher in participants with myeloid
336 CH (53% smokers) than those without myeloid CH (44% smokers; $P_{FDR}=3.38 \times 10^{-6}$, Table 5,
337 Supplementary Table 7). This effect was associated with current smoking status (OR=1.10,
338 $P_{FDR}=6.14 \times 10^{-6}$) rather than past smoking status (OR=1.02, $P_{FDR}=0.08$).

339

340 Strikingly, breakdown of myeloid CH by specific mutation type revealed that variants in
341 *ASXL1* are strongly associated with current smoking status (OR=1.07, $P_{FDR}=1.92 \times 10^{-5}$) and the
342 only abnormality associated with past smoking status (OR=1.04, $P_{FDR}=0.0026$). Indeed, 69%
343 of participants with *ASXL1* mutations were past or current smokers. Participants with
344 myeloid CH without *ASXL1* mutations (n=1066) remained significantly associated with
345 current smoking but the effect was weaker (OR=1.07, $P_{FDR}=0.0008$). Both *TET2* and *DNMT3A*
346 variants showed a significant, but relatively modest, association with current smoking status
347 but there was no discernible association between smoking and variants in *JAK2*, *SRSF2*,
348 *PPM1D* or for acquired myeloid mCA (Table 5; Supplementary Table 7; all comparisons were
349 corrected for gender and age).

350

351 *The association between myeloid CH and clinical phenotype, blood traits and biochemical*
352 *measures*

353 We sought to determine if myeloid CH in hematological malignancy-free participants
354 (n=911) was associated with 29 blood features, 30 biochemical markers and 395 non-
355 malignant clinical phenotypes. Myeloid CH was associated with 9 blood features including
356 red cell distribution width (RDW), estimated to be 1.02 fold higher in myeloid CH cases
357 ($P_{FDR}=9.7 \times 10^{-4}$; linear regression; Table 6), in line with previous findings.⁸ In addition,
358 myeloid CH was associated with a significant increase in all platelet indices, and particularly
359 platelet distribution width (PDW; OR=1.03, $P_{FDR}=6 \times 10^{-6}$) and a decrease in basophil indices,
360 most notably basophil counts (OR=0.95, $P_{FDR}=5.9 \times 10^{-4}$, Table 6). Myeloid CH cases also
361 showed an anemia-like blood profile with low hemoglobin and a decrease in cholesterol,
362 high density lipoprotein and creatinine (Table 7). Clinically, a significant association was seen
363 with urinary tract disorders defined by ICD-10 codes N35.9 (urethral stricture) and N32.0
364 (bladder constriction) (Table 8).

365

366 Looking at individual drivers of myeloid CH, several associations emerge. The strongest
367 multifactorial findings were for *JAK2 V617F*, which was associated with increased red cell
368 and platelet parameters (Table 6). *ASXL1* mutations were associated with increased RDW
369 and PDW, an anemia-like profile and reduced IGF-1. *TET2* mutations were associated with
370 reduced eosinophils, increased monocytes and, clinically, ulcer of the lower limb,
371 agranulocytosis and chronic obstructive pulmonary disease. Participants with *SRSF2* P95
372 showed an increase in the percentage of reticulocytes and a decrease in HDL cholesterol,
373 and mCA was associated with multiple parameters, particularly low basophils, increased
374 PDW, several decreased biochemical measures as well as urinary tract disorders. Finally,
375 *DNMT3A* and *PPM1D* mutations did not show any associations except for a relatively minor
376 increase in platelet and monocyte counts, respectively. The complete results are shown in
377 Supplementary Tables 8, 9 and 10.

378

379 We noted that the incidence of all-cause mortality since study entry was higher in
380 participants with myeloid CH without hematological malignancy (42/911; 4.6%) compared
381 with controls (674/30892; 2.2%). Using multivariable Cox regression, myeloid CH was shown
382 to be associated with an increased risk of all-cause mortality (HR=1.44; $P=0.021$; mean

383 follow up=8.1 years) after adjusting for age and sex. There was no association between CH
384 and subsequent MI or stroke in the absence of a hematological malignancy (MI: myeloid CH,
385 n=873 vs controls, n=30271; HR=1.16; P=0.53; stroke: myeloid CH, n=890 vs controls,
386 n=30472; HR=1.18; P=0.58).

387

388

389 **DISCUSSION**

390

391 The prevalence and significance of CH has been reported in several cohorts, but our study
392 has several distinctive features. UK Biobank is a very large population-based cohort that
393 includes an extensive repertoire of baseline phenotypic data as well as >9 year prospective
394 collection of clinical follow up information. Genome wide SNP data is available for the great
395 majority of participants (n=486,941), and WES data for a subset (n=49,956), all derived from
396 a single baseline peripheral sample taken at study entry. Thus, we were able to assess CH
397 associated with both mCA and somatic mutations, albeit with a modest limit of detection
398 compared to some published studies. We focused on myeloid mCA and genes known to be
399 mutated in myeloid disorders with the specific aim of understanding the causes and
400 consequences of myeloid CH.

401

402 In total we identified 1,166 subjects with myeloid CH. Of these, 678 had mutations in one or
403 more genes (1.4% of cases who underwent WES), and 506 had myeloid mCA (0.1% of
404 subjects who had a SNP array; 18 subjects had both mCA and somatic mutations). Given the
405 median age of the UK Biobank cohort (58 years), the limited sensitivity afforded by WES and
406 the fact we focused only on myeloid mCA, the prevalence of CH is broadly comparable to
407 published reports.^{6,8}

408

409 Like other studies, we found that *DNMT3A* was the most commonly mutated gene and that
410 there was a strong relationship between both mCA and mutations with age. Clinically, we
411 confirmed associations between CH and all-cause mortality and a specific association
412 between *TET2* mutations and chronic obstructive pulmonary disorders with acute lower

413 respiratory infection.¹³ In common with some other studies,^{7, 13} we did not find an
414 association between CH and MI or stroke. The reason that the association between CH and
415 cardiovascular disease is very prominent in some studies¹² but not others is presumably
416 explained by differences in cohort structure, follow up time, and definitions of CH. UK
417 Biobank had an upper recruitment age of 69 years and the follow up is only 9.1 years. Our
418 analysis is estimated to have 86% power (Supplementary Figure 5) to detect an association
419 between CH and MI based on an HR of 1.9, as previously reported,¹² and an overall event
420 rate of 1.4% (439/31144). Our definition of CH included both chromosomal and mutational
421 events, with a stringent definition of pathogenicity for mutations, and all abnormalities
422 being present at a clonal fraction >10%. Clearly, as more UK Biobank cases are sequenced
423 and the median follow-up is extended, more associations are likely to emerge.

424

425 We also confirmed an association with CH and smoking⁷ and showed for the first time that
426 this effect is predominantly, but not exclusively, associated with *ASXL1* mutations. Whilst
427 this association, and other novel findings in our study, needs to be confirmed in an
428 independent population-based cohort, we note that *ASXL1* mutations have recently been
429 associated with smoking in a large cohort of post-therapy cancer patients.²⁸ CH has
430 previously been associated chemotherapy and radiotherapy^{29, 30} which, along with the
431 association with ageing, suggest that a link between stress hematopoiesis and the
432 development of clonality. Although it is conceivable that smoking preferentially induces
433 *ASXL1* mutations, it seems more likely that the chronic inflammation induced by smoking
434 creates an environment within which *ASXL1* mutant clones have a selective advantage. A
435 previous study noted an increase in the incidence of C>A transversions in smokers³⁰ but we
436 found the C>A transversion rate in *ASXL1* was similar in smokers (18%) compared to non-
437 smokers (17%). Overall, C>T transitions (n=245; 66%) represented the most common single
438 nucleotide substitution, as expected.³¹ Of note, we found that CH was associated with
439 inflammatory related disorders (e.g. chronic obstructive pulmonary disorders,
440 agranulocytosis, and ulcer of lower limb) but it is unclear whether these might be caused by
441 CH or whether the association reflects a common inflammatory background.

442

443 Whilst extrinsic factors play a role in the development of CH, constitutional genetics is also
444 important. Genome wide association analysis identified rs34002450, an intronic variant in
445 *TERT* gene, to be associated with CH in the Icelandic population.⁷ This SNP was not included
446 on the array platform used by UK Biobank, but we independently identified two distinct
447 signals within *TERT* that achieved genome wide significance: rs7726159 and rs2853677. One
448 of these SNPs, rs7726159, is in LD with rs34002450 ($r^2= 0.70$) whereas rs2853677 has been
449 associated with MPN and *JAK2* V617F associated CH.³² Two additional SNPs in *TERT* have
450 been reported to predispose to the development of MPN and *JAK2* V617F associated CH:
451 rs2736100^{33, 34} and rs7705526.³² LD analysis for the two primary signals (rs7726159 and
452 rs2853677; Supplementary Table 12) revealed (i) rs7726159 is in LD with rs7705526
453 ($r^2=0.79$) but does not reach genome-wide significance for association with self-reported
454 polycythemia vera (PV) in UK Biobank ($P=2.5 \times 10^{-4}$; <http://big.stats.ox.ac.uk>); (ii) rs2853677 is
455 not in LD with rs7705526, ($r^2=0.19$), but is associated with PV ($P=2.9 \times 10^{-6}$;
456 <http://big.stats.ox.ac.uk>); (iii) rs2736100 is in modest LD with rs7705526 ($r^2=0.51$) and is
457 associated with PV ($P=1 \times 10^{-5}$; <http://big.stats.ox.ac.uk>). Thus, it appears that variation in
458 intron 2 (rs7726159) is associated with CH but does not predict development of MPN but
459 variation in intron 3 (rs2736100, rs2853677 and rs7705526) does predict development of
460 MPN (Figure 2; Supplementary Table 12). SNP rs2853677 is not in LD with any of the other
461 variants and is thus a unique independent signal for both CH and MPN.

462

463 We conclude that both genetic and environmental factors play an important role in the
464 development of CH. Smoking history is strongly associated with *ASXL1* mutated CH and
465 genetic variation at *TERT* may predispose to CH independently of predisposition to MPN.
466 *TERT* encodes telomerase reverse transcriptase and is essential for telomere maintenance,
467 but it also appears to function as a transcriptional co-activator³⁵ and impacts on the tumor
468 microenvironment via diverse pathways, including inflammation. It is possible therefore that
469 chronic inflammation provides a link between genetic and environmental predisposition to
470 CH.

471

472 **ACKNOWLEDGEMENTS**

473 AAZD was supported by Lady Tata International Award; NCPC was supported by Blood

474 Cancer UK

475

476 **DISCLOSURE OF CONFLICTS OF INTEREST**

477 The authors have no conflicts of interest to declare

478 **REFERENCES**

479

- 480 1. Busque L, Mio R, Mattioli J, Brais E, Blais N, Lalonde Y, *et al.* Nonrandom X-inactivation
481 patterns in normal females: lyonization ratios vary with age. *Blood* 1996; **88**(1): 59-65.
- 482 2. Busque L, Patel JP, Figueroa ME, Vasanthakumar A, Provost S, Hamilou Z, *et al.* Recurrent
483 somatic TET2 mutations in normal elderly individuals with clonal hematopoiesis. *Nature*
484 *Genetics* 2012; **44**(11): 1179.
- 486 3. Xu X, Zhang Q, Luo J, Xing S, Li Q, Krantz SB, *et al.* JAK2V617F: prevalence in a large Chinese
487 hospital population. *Blood* 2007; **109**(1): 339-342.
- 489 4. Laurie CC, Laurie CA, Rice K, Doheny KF, Zelnick LR, McHugh CP, *et al.* Detectable clonal
490 mosaicism from birth to old age and its relationship to cancer. *Nature Genetics* 2012; **44**(6):
491 642.
- 493 5. Jacobs KB, Yeager M, Zhou W, Wacholder S, Wang Z, Rodriguez-Santiago B, *et al.* Detectable
494 clonal mosaicism and its relationship to aging and cancer. *Nature Genetics* 2012; **44**(6): 651.
- 496 6. Genovese G, Kähler AK, Handsaker RE, Lindberg J, Rose SA, Bakhoum SF, *et al.* Clonal
497 hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *New England*
498 *Journal of Medicine* 2014; **371**(26): 2477-2487.
- 500 7. Zink F, Stacey SN, Norddahl GL, Frigge ML, Magnusson OT, Jonsdottir I, *et al.* Clonal
501 hematopoiesis, with and without candidate driver mutations, is common in the elderly.
502 *Blood* 2017; **130**(6): 742-752.
- 504 8. Jaiswal S, Fontanillas P, Flannick J, Manning A, Grauman PV, Mar BG, *et al.* Age-related clonal
505 hematopoiesis associated with adverse outcomes. *New England Journal of Medicine* 2014;
506 **371**(26): 2488-2498.
- 508 9. Acuna-Hidalgo R, Sengul H, Steehouwer M, van de Vorst M, Vermeulen SH, Kiemeny LA, *et al.*
509 Ultra-sensitive sequencing identifies high prevalence of clonal hematopoiesis-associated
510 mutations throughout adult life. *The American Journal of Human Genetics* 2017; **101**(1): 50-
511 64.
- 513 10. Young AL, Challen GA, Birman BM, Druley TE. Clonal haematopoiesis harbouring AML-
514 associated mutations is ubiquitous in healthy adults. *Nature Communications* 2016; **7**:
515 12484.
- 516

517

- 518 11. Steensma DP, Bejar R, Jaiswal S, Lindsley RC, Sekeres MA, Hasserjian RP, *et al.* Clonal
519 hematopoiesis of indeterminate potential and its distinction from myelodysplastic
520 syndromes. *Blood* 2015; **126**(1): 9-16.
- 521
522 12. Jaiswal S, Natarajan P, Silver AJ, Gibson CJ, Bick AG, Shvartz E, *et al.* Clonal hematopoiesis
523 and risk of atherosclerotic cardiovascular disease. *New England Journal of Medicine* 2017;
524 **377**(2): 111-121.
- 525
526 13. Buscarlet M, Provost S, Zada YF, Barhdadi A, Bourgoin V, Lépine G, *et al.* DNMT3A and TET2
527 dominate clonal hematopoiesis and demonstrate benign phenotypes and different genetic
528 predispositions. *Blood* 2017; **130**(6): 753-762.
- 529
530 14. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, *et al.* The UK Biobank resource
531 with deep phenotyping and genomic data. *Nature* 2018; **562**(7726): 203.
- 532
533 15. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, *et al.* Genome-wide genetic
534 data on ~ 500,000 UK Biobank participants. *BioRxiv* 2017: 166298.
- 535
536 16. Staaf J, Lindgren D, Vallon-Christersson J, Isaksson A, Göransson H, Juliusson G, *et al.*
537 Segmentation-based detection of allelic imbalance and loss-of-heterozygosity in cancer cells
538 using whole genome SNP arrays. *Genome Biology* 2008; **9**(9): R136.
- 539
540 17. Van Hout CV, Tachmazidou I, Backman JD, Hoffman JX, Yi B, Pandey A, *et al.* Whole exome
541 sequencing and characterization of coding variation in 49,960 individuals in the UK Biobank.
542 *bioRxiv* 2019: 572347.
- 543
544 18. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, *et al.* The sequence
545 alignment/map format and SAMtools. *Bioinformatics* 2009; **25**(16): 2078-2079.
- 546
547 19. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-
548 throughput sequencing data. *Nucleic Acids Research* 2010; **38**(16): e164-e164.
- 549
550 20. Hollander M, Wolfe DA, Chicken E. *Nonparametric statistical methods*, vol. 751. John Wiley
551 & Sons, 2013.
- 552
553 21. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, *et al.* PLINK: a tool set
554 for whole-genome association and population-based linkage analyses. *The American Journal*
555 *of Human Genetics* 2007; **81**(3): 559-575.
- 556
557 22. Millard LA, Davies NM, Gaunt TR, Davey Smith G, Tilling K. Software Application Profile:
558 PHESANT: a tool for performing automated phenome scans in UK Biobank. *International*
559 *Journal of Epidemiology* 2017; **47**(41): 29–35.

560
561 23. Therneau TM, Grambsch PM. The Cox model. *Modeling survival data: extending the Cox*
562 *model*. Springer 2000, pp 39-77.

563
564 24. Grolemund G, Wickham H. Dates and times made easy with lubridate. *Journal of Statistical*
565 *Software* 2011; **40**(3): 1-25.

566
567 25. Roman E, Smith A, Appleton S, Crouch S, Kelly R, Kinsey S, *et al*. Myeloid malignancies in the
568 real-world: Occurrence, progression and survival in the UK's population-based
569 Haematological Malignancy Research Network 2004–15. *Cancer Epidemiology* 2016; **42**: 186-
570 198.

571
572 26. McKerrell T, Park N, Moreno T, Grove CS, Ponstingl H, Stephens J, *et al*. Leukemia-associated
573 somatic mutations drive distinct patterns of age-related clonal hemopoiesis. *Cell Reports*
574 2015; **10**(8): 1239-1245.

575
576 27. Jones AV, Chase A, Silver RT, Oscier D, Zoi K, Wang YL, *et al*. JAK2 haplotype is a major risk
577 factor for the development of myeloproliferative neoplasms. *Nature Genetics* 2009; **41**(4):
578 446-449.

579
580 28. Bolton KL, Ptashkin RN, Gao T, Braunstein L, Devlin SM, Kelly D, *et al*. Oncologic therapy
581 shapes the fitness landscape of clonal hematopoiesis. *bioRxiv* 2019: 848739.

582
583 29. Gillis NK, Ball M, Zhang Q, Ma Z, Zhao Y, Yoder SJ, *et al*. Clonal haemopoiesis and therapy-
584 related myeloid malignancies in elderly patients: a proof-of-concept, case-control study. *The*
585 *Lancet Oncology* 2017; **18**(1): 112-121.

586
587 30. Coombs CC, Zehir A, Devlin SM, Kishtagari A, Syed A, Jonsson P, *et al*. Therapy-related clonal
588 hematopoiesis in patients with non-hematologic cancers is common and associated with
589 adverse clinical outcomes. *Cell Stem Cell* 2017; **21**(3): 374-382. e374.

590
591 31. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SA, Behjati S, Biankin AV, *et al*. Signatures of
592 mutational processes in human cancer. *Nature* 2013; **500**(7463): 415-421.

593
594 32. Hinds DA, Barnholt KE, Mesa RA, Kiefer AK, Do CB, Eriksson N, *et al*. Germ line variants
595 predispose to both JAK2 V617F clonal hematopoiesis and myeloproliferative neoplasms.
596 *Blood, The Journal of the American Society of Hematology* 2016; **128**(8): 1121-1128.

597
598 33. Oddsson A, Kristinsson S, Helgason H, Gudbjartsson D, Masson G, Sigurdsson A, *et al*. The
599 germline sequence variant rs2736100_C in TERT associates with myeloproliferative
600 neoplasms. *Leukemia* 2014; **28**(6): 1371-1374.

601

602 34. Tapper W, Jones AV, Kralovics R, Harutyunyan AS, Zoi K, Leung W, *et al.* Genetic variation at
603 MECOM, TERT, JAK2 and HBS1L-MYB predisposes to myeloproliferative neoplasms. *Nature*
604 *Communications* 2015; **6**(1): 1-11.

605
606 35. Liu N, Guo XH, Liu JP, Cong YS. Role of telomerase in the tumour microenvironment. *Clinical*
607 *and Experimental Pharmacology and Physiology* 2019; **47**(43):357-364.

608

609

610

611 **FIGURE LEGENDS**

612

613 Figure 1: The relationship between clonal hematopoiesis and age. (A) Total mCA frequency
614 across different age intervals showing increased risk of acquiring a mCA by 1.02 fold
615 ($P=1.80 \times 10^{-19}$). (B) box plot showing increased age in subjects with ≥ 1 mCA (median = 60
616 years; $n=5,040$) compared to those with no mCA (median = 58 years; $n=481,901$; $P=1.80 \times 10^{-19}$). (C) Myeloid mCA frequency across different age intervals showing increased risk of
617 acquiring a myeloid mCA by 1.1 fold ($P=1.57 \times 10^{-38}$). (D) box plot showing increased age in
618 subjects with ≥ 1 myeloid mCA (median = 63 years; $n=506$) compared to those with no mCA
619 (median = 58 years; $n=481,901$; $P=1.57 \times 10^{-38}$). (E) Frequency of individual mutations
620 showing an age-related increase for all genes and increased risk of acquiring mutations by
621 1.1 fold ($P=5.89 \times 10^{-47}$). (F) box plot showing increased age in subjects with ≥ 1 mutation
622 (median = 63 years; $n=678$) compared to those with no mutations (median = 58 years;
623 $n=49,278$; $P=5.89 \times 10^{-47}$).

624

625
626 Figure 2: Genome wide association analysis to identify genetic predisposition to myeloid
627 clonal hematopoiesis. (A) Quantile-quantile plot showing observed versus expected P
628 values. No evidence was seen for systematic bias between cases and controls, or population
629 stratification ($\lambda=1.021$) (B) Manhattan plot summarising the significance of SNPs
630 across the genome. The red line indicates genome wide significance ($P < 5 \times 10^{-8}$) and the blue
631 line indicates values that were of suggestive significance ($P < 10^{-5}$). Clusters related to *TERT*
632 and *JAK2* are indicated. (C) Locus zoom plot focusing on SNPs in the region of *TERT* at
633 chromosome band 5p15. (D) Positions and genetic relationships between *TERT* SNPs.

Table 1: Summary of UK Biobank cohort.

	Males n (%)	Females n (%)	Total
Phenotypic data ¹	229,129 (46)	273,395 (54)	502,524
Myeloid disorders ²	1,157 (57)	873 (43)	2,030
Lymphoid disorders ³	5,747 (44)	7,390 (56)	13,137
Other cancers	49,435 (41)	71,420 (59)	120,855
Cancer free	172,790 (47)	193,712 (53)	366,502
SNP array data ⁴	222,858 (46)	264,083 (54)	486,941
Myeloid disorders	1,097 (57)	816 (43)	1,913
Lymphoid disorders	5,566 (44)	6,980 (56)	12,546
Other cancers	48,101 (41)	68,820 (59)	116,921
Cancer free	168,094 (47)	187,467 (53)	355,561
WES data ⁵	22,714 (45)	27,242 (55)	49,956
Myeloid disorders	97 (53)	85 (47)	182
Lymphoid disorders	473 (46)	550 (54)	1,023
Other cancers	4,972 (41)	7,265 (59)	12,237
Cancer free	17,172 (47)	19,342 (53)	36,514

- 1) Includes cases who had the specified disorder at any time during the study period.
- 2) Of 2030 participants with a myeloid disorder, 315 were also diagnosed with another non-myeloid haematological disorder during the study period.
- 3) 34 cases with unspecified haematological malignancy were included in the lymphoid group
- 4) Data available from 488,377 cases of which 1436 were excluded following QC (see Methods).
- 5) Data available from 49,996 cases of which 40 were excluded following QC (see Methods).

Table 2: Summary of mCA identified across the cohort.

Group	SNP array samples	Total mCA events (per sample)	P*	Samples with at least one event											
				mCA			aUPD			CNG			CNL		
				n (%)	OR	P _{FDR}	n	OR	P _{FDR}	n	OR	P _{FDR}	n	OR	P _{FDR}
Myeloid disorders	1913	316 (1.5)	4.10x10 ⁻⁶	210 (11)	13.21	3.74x10 ⁻¹⁴⁵	158	16.39	8.78x10 ⁻¹²⁴	37	40.88	3.10x10 ⁻⁴⁴	34	4.61	2.18x10 ⁻¹²
lymphoid disorders	12546	768 (2.1)	0.54	363 (2.9)	3.19	1.41x10 ⁻¹⁰⁵	146	2.15	1.77x10 ⁻¹⁹	81	14.00	3.09x10 ⁻⁵⁵	194	4.01	1.22x10 ⁻⁸³
Other cancers	116921	1854 (1.6)	0.89	1185 (1)	1.10	4.00x10 ⁻³	657	1.03	0.26	67	1.24	0.09	527	1.16	3.6x10 ⁻³
Cancer free	355561	5269 (1.6)		3282 (0.9)			1938			165			1386		

The number of mosaic chromosomal abnormalities (mCA) identified in each phenotypic group out of the total number of samples with SNP array data passing QC. The number of events for each mCA subcategory are also shown; acquired uniparental disomy (aUPD), copy number gain (CNG) and copy number loss (CNL). The mean number of mCA events in participants with either myeloid, lymphoid, or other cancers were compared with cancer free controls using Mann Whitney U tests (P*). Fisher's exact tests were used to compare the number of events which were corrected for 12 tests using the false discovery rate (P_{FDR}).

Table 3: Summary of mCA events significantly associated with myeloid disorders

Event	mCA type	Cancer free, No. positive [†]	Myeloid malignancies*		
			No. positive	OR	P _{FDR}
chr9p	aUPD	7	86	2859	6.30 x10 ⁻¹⁹¹
chr9p	CNG	1	15	3335	7.93 x10 ⁻³³
chr14q	aUPD	23	12	116	2.67 x10 ⁻¹⁸
chr9q	CNG	1	8	1771	6.93 x10 ⁻¹⁷
chr1p	aUPD	31	10	72	1.25 x10 ⁻¹³
chr20i	CNL	36	10	63	3.31 x10 ⁻¹³
chr4q	aUPD	13	8	136	1.04 x10 ⁻¹²
chr1q	CNG	1	4	833	4.26 x10 ⁻⁸
chr8p	CNG	5	4	177	8.45 x10 ⁻⁷
chr9p	CNL	1	3	662	5.82 x10 ⁻⁶
chr8q	CNG	4	3	166	4.03 x10 ⁻⁵
chr7q	aUPD	6	3	110	1.01 x10 ⁻⁴
chr7q	CNL	0	2	-	2.31 x10 ⁻⁴
chr17p	aUPD	4	2	110	2.85 x10 ⁻³
chr19q	aUPD	12	2	37	0.01
chr11q	aUPD	24	2	18	0.04
chr22q	aUPD	27	2	16	0.05

[†]From a total of 355,561 cancer free sample

*From a total of 1,614 samples with a myeloid malignancy. 299/1913 myeloid cases were excluded from this analysis because they had both myeloid and lymphoid disorders

Table 4: Summary of somatic mutations.

	Mutations		Participants				
	N	VAF median (range)	Total	Myeloid n= 182	Lymphoid n=1,023	Other Cancer n=12,237	Cancer Free n=36,514
<i>DNMT3A</i> LOF	223	0.17 (0.07-0.50)	222	1	5	64	152
<i>DNMT3A</i> R882	86	0.17 (0.11-0.40)	86	1	1	25	59
<i>TET2</i> LOF	223	0.18 (0.06-0.68)	208	9	10	55	134
<i>ASXL1</i> LOF	101	0.21 (0.08-0.49)	100	4	3	24	69
<i>JAK2</i> V617F	40	0.27 (0.12-0.90)	40	25	0	4	11
<i>SRSF2</i> P95	20	0.24 (0.11-0.47)	20	5	0	2	13
<i>PPM1D</i> LOF	28	0.21 (0.10-0.51)	28	0	2	8	18
TOTAL	721	0.19 (0.08-0.90)	678	41	21	174	442

Table 5: The relationship between smoking and clonal hematopoiesis

Marker ¹	No. myeloid	No. of smokers ²		Previous smoking ³ “Data-Field 1249”		Current smoking ³ “Data-Field 1239”		Combined smoking ⁴ “Data-Field 20116”	
		Past	Current	OR	P_{FDR}	OR	P_{FDR}	OR ^a ; OR ^b	P_{FDR}
mCA	506	218	48	1.01	0.39	1.03	0.18	1.19; 1.42	0.091
<i>ASXL1</i> LOF	100	49	20	1.04	2.60×10^{-3}	1.07	1.92×10^{-5}	1.94; 4.68	1.02×10^{-5}
<i>DNMT3A</i> LOF or R882	308	117	35	1.00	1.00	1.05	0.03	1.03; 1.64	0.07
<i>JAK2 V617F</i> or chr9p mCA	155	64	7	0.99	0.68	0.95	0.18	0.95; 0.54	0.27
<i>PPM1D</i> LOF	28	15	3	1.02	0.16	1.01	0.68	2.07; 2.05	0.23
<i>SRSF2</i> P95	20	10	1	1.01	0.55	1.00	1.00	0.88; 1.82	0.83
<i>TET2</i> LOF	208	75	27	0.99	0.5	1.06	6.40×10^{-3}	0.88; 1.82	0.03
All myeloid CH	1166	488	134	1.02	0.09	1.10	6.14×10^{-6}	1.17; 1.76	3.38×10^{-6}
Myeloid CH without <i>ASXL1</i>	1066	439	114	1.01	0.36	1.07	8.8×10^{-4}	1.12; 1.59	5.81×10^{-4}

- 1) Loss of function (LOF); clonal hematopoiesis (CH); mosaic chromosomal alterations (mCA)
- 2) Number of smokers encoded in the combined smoking status in UK Biobank “Data-Field 20116”
- 3) Results of ordinal logistic regression, total tests = 16, corrected for age, sex and FDR.
- 4) Results of multinomial logistic regression, total tests = 8, corrected for age, sex and FDR. Odds ratios are estimated for past smoking level (a), and current smoking (b)

Table 6: Blood features associated with myeloid markers*

Group	Blood feature	Units	No.	Mean value		OR (CI 97.5%)	P _{FDR}
				Cases	Control		
mCA	Basophil count	10 ⁹ /L	30272	0.03	0.04	0.92 (0.90-0.94)	9.1x10 ⁻¹³
	Platelet distribution width	%	30282	16.69	16.46	1.04 (1.03-1.05)	4.3x10 ⁻⁹
	Hematocrit percentage	%	30282	41.05	41.61	0.98 (0.97-0.99)	1.0x10 ⁻⁴
	Basophil percentage	%	30272	0.53	0.61	0.97 (0.96-0.99)	2.9x10 ⁻⁴
	Mean corpuscular hemoglobin	g/dL	30282	34.43	34.27	1.02 (1.01-1.03)	4.2x10 ⁻³
	Red blood cell count	10 ⁹ /L	30282	4.47	4.54	0.98 (0.97-0.99)	4.4x10 ⁻³
	Red blood cell distribution width	%	30282	13.81	13.5	1.02 (1.01-1.03)	4.9x10 ⁻³
	Hemoglobin concentration	g/dL	30282	14.13	14.26	0.98 (0.98-0.99)	6.8x10 ⁻³
	Mean sphered cell volume	fL	28916	83.99	84.5	0.98 (0.97-0.99)	3.0x10 ⁻²
ASXL1	Platelet distribution width	%	30082	16.72	16.46	1.03 (1.01-1.04)	5.9x10 ⁻⁴
	Red blood cell distribution width	%	30082	13.93	13.5	1.03 (1.01-1.04)	7.0x10 ⁻⁴
	Mean corpuscular volume	fL	30082	90.68	91.8	0.98 (0.97-0.99)	9.1x10 ⁻⁴
	Mean corpuscular haemoglobin	pg	30082	31.03	31.47	0.98 (0.97-0.99)	2.8x10 ⁻³
	Mean sphered cell volume	fL	28712	83.36	84.5	0.98 (0.97-0.99)	1.0x10 ⁻²
DNMT3A	Platelet count	10 ⁹ /L	30289	251.27	242.76	1.02 (1.01-1.03)	1.8x10 ⁻²
JAK2	Platelet crit	%	30019	0.30	0.22	1.04 (1.03-1.06)	7.5x10 ⁻¹²
	Platelet count	10 ⁹ /L	30019	341.46	242.76	1.04 (1.03-1.06)	1.3x10 ⁻¹¹
	Red blood cell distribution width	%	30019	15.31	13.5	1.04 (1.03-1.05)	3.6x10 ⁻⁹
	Platelet distribution width	%	30019	17.14	16.46	1.03 (1.02-1.05)	1.0x10 ⁻⁶
	High light scatter reticulocyte	10 ¹² /L	28656	0.02	0.02	1.02 (1.01-1.03)	4.0x10 ⁻²
PPM1D	Monocyte count	10 ⁹ /L	30007	0.59	0.48	1.02 (1.01-1.03)	3.5x10 ⁻²
SRSF2	Reticulocyte percentage	%	28643	1.88	1.32	1.02 (1.01-1.03)	1.3x10 ⁻²
	High light scatter reticulocyte	%	28643	0.62	0.4	1.02 (1.01-1.03)	2.4x10 ⁻²
TET2	Eosinophil count	10 ⁹ /L	30169	0.15	0.17	0.98 (0.97-0.99)	1.1x10 ⁻³
	Eosinophil percentage	%	30169	2.18	2.53	0.98 (0.97-0.99)	4.4x10 ⁻³
	Monocyte percentage	%	30169	7.83	7.06	1.02 (1.01-1.03)	3.0x10 ⁻²
All Myeloid CH	Platelet distribution width	%	30883	16.57	16.46	1.03 (1.02-1.04)	6.0x10 ⁻⁶
	Basophil count	10 ⁹ /L	30873	0.04	0.04	0.95 (0.93-0.97)	5.9x10 ⁻⁴
	Red blood cell distribution width	%	30883	13.71	13.5	1.02 (1.01-1.04)	9.7x10 ⁻⁴
	Platelet crit	%	30883	0.23	0.22	1.02 (1.01-1.03)	1.1x10 ⁻³
	Hematocrit percentage	%	30883	41.46	41.61	0.98 (0.97-0.99)	1.7x10 ⁻³
	Platelet count	10 ⁹ /L	30883	248.32	242.76	1.02 (1.01-1.03)	3.3x10 ⁻³
	Hemoglobin concentration	g/dL	30883	14.22	14.26	0.99 (0.98-0.99)	1.7x10 ⁻²
	Basophil percentage	%	30883	0.58	0.61	0.98 (0.97-0.99)	4.4x10 ⁻²

* Total number of linear regression tests =232, corrected for age, sex, and smoking status.

Table 7: Biochemical measures associated with myeloid markers*

Group	Biochemistry measure	Units	N	Mean in		OR (CI 97.5%)	P_{FDR}
				Cases	Control		
mCA	Creatinine	μmol/L	29280	71.452	72.686	0.98 (0.97-0.99)	0.001
	Apolipoprotein A	g/L	27335	1.517	1.555	0.98 (0.97-0.99)	0.004
	Phosphate	mmol/L	27515	1.169	1.200	0.98 (0.97-0.99)	0.005
	HDL cholesterol	mmol/L	27546	1.420	1.474	0.98 (0.97-0.99)	0.010
	Albumin	g/L	27576	44.848	45.518	0.98 (0.97-0.99)	0.018
ASXL1	IGF-1	nmol/L	28977	19.043	21.697	0.98 (0.97-0.99)	0.033
SRSF2	HDL cholesterol	mmol/L	27294	1.237	1.474	0.98 (0.97-0.99)	0.027
All CH	Cholesterol	mmol/L	29874	5.619	5.697	0.98 (0.97-0.99)	0.033
	HDL cholesterol	mmol/L	28085	1.450	1.474	0.98 (0.97-0.99)	0.040
	Creatinine	μmol/L	29851	72.706	72.686	0.99 (0.98-0.99)	0.041

* Total number of linear regression tests =150, corrected for age, sex, and smoking status.

Table 8: Clinical phenotypes associated with myeloid markers

Marker	Phenotype	Cases	Controls n=30,892	OR (CI 97.5%)	P_{FDR}
mCA n=301	N35.9 Urethral stricture, unspecified	10 (3.3%)	189 (0.6%)	1.17 (1.09-1.24)	0.008
	N32.0 Bladder-neck obstruction	7 (2.3%)	98 (0.03%)	1.19 (1.09-1.28)	0.009
TET2 n=189	L97 Ulcer of lower limb, not elsewhere classified	4 (2.1%)	19 (0.06%)	1.28 (1.14-1.39)	0.004
	D70 Agranulocytosis	4 (2.1%)	51 (0.16%)	1.23 (1.10-1.34)	0.009
	J44.0 COPD with acute lower respiratory infection	7 (3.7%)	137 (0.44%)	1.16 (1.07-1.23)	0.009

* Total number of logistic regression tests =3160, corrected for age, sex, and smoking status.

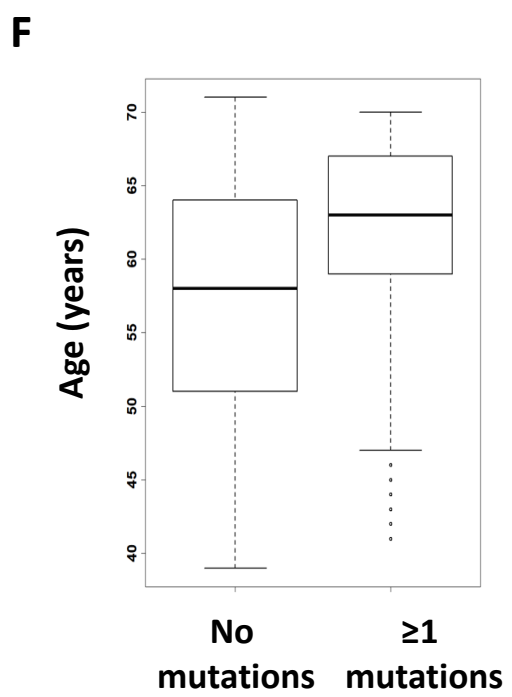
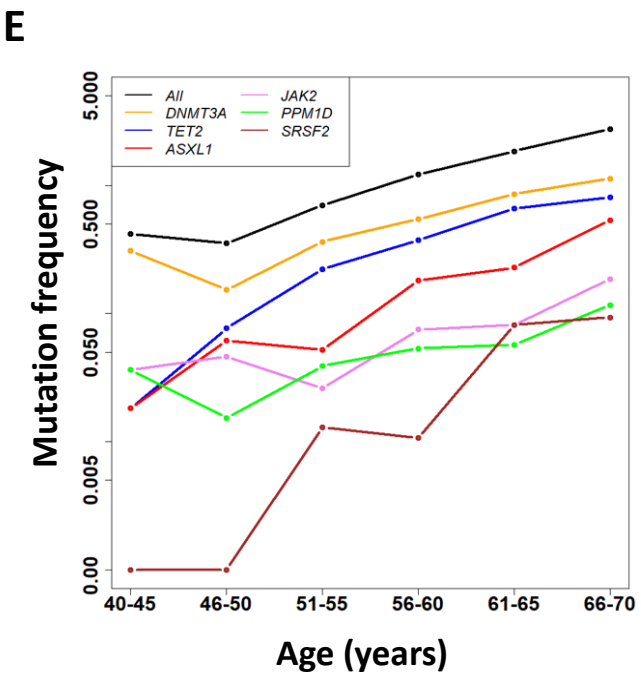
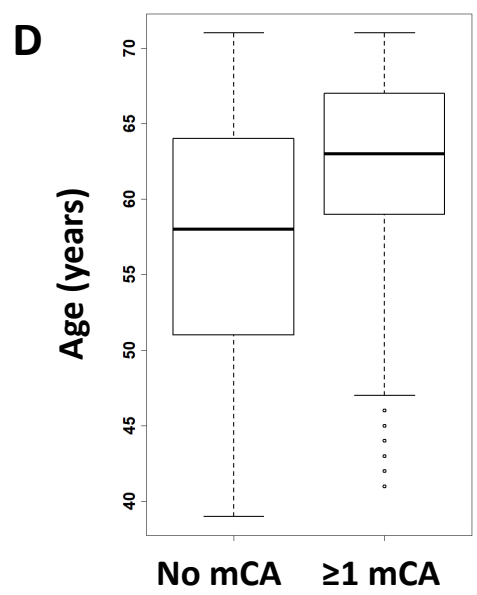
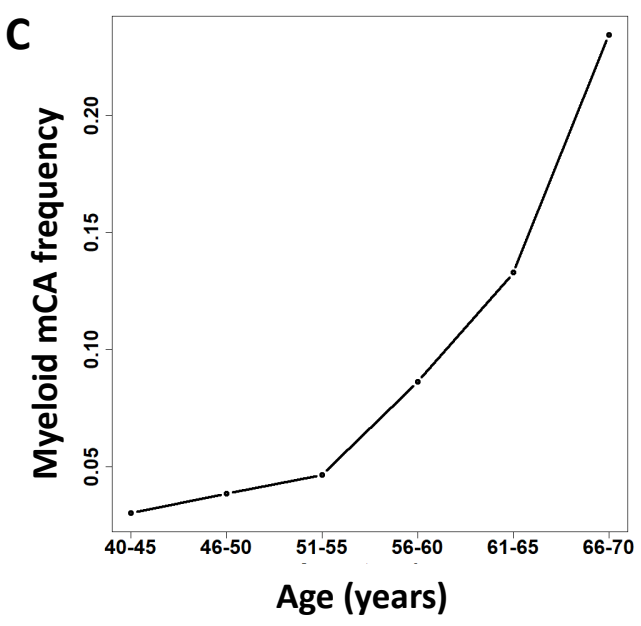
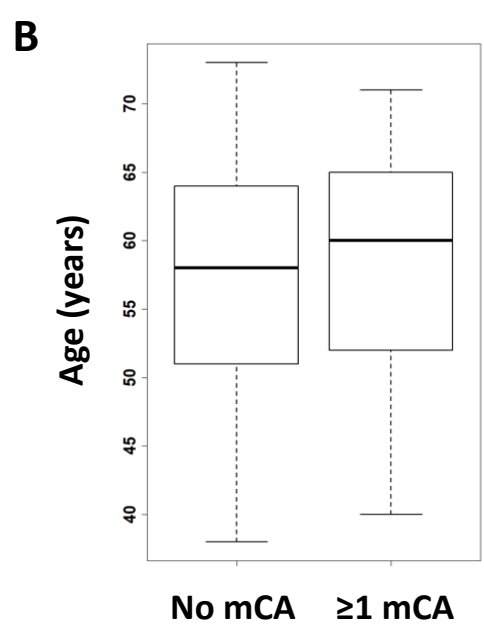
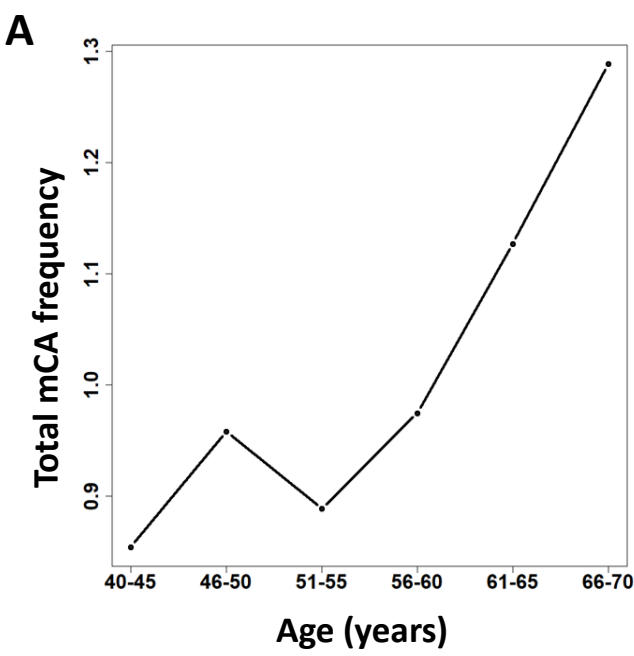
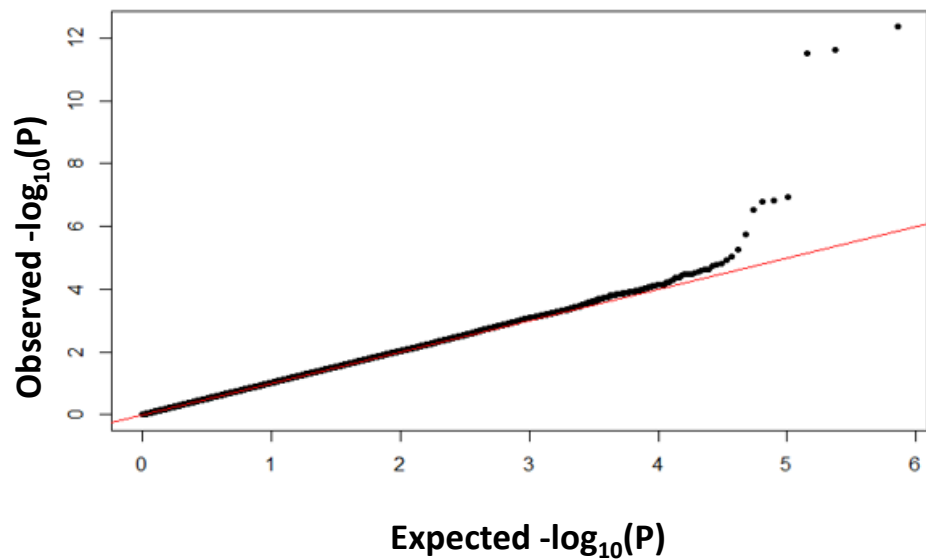
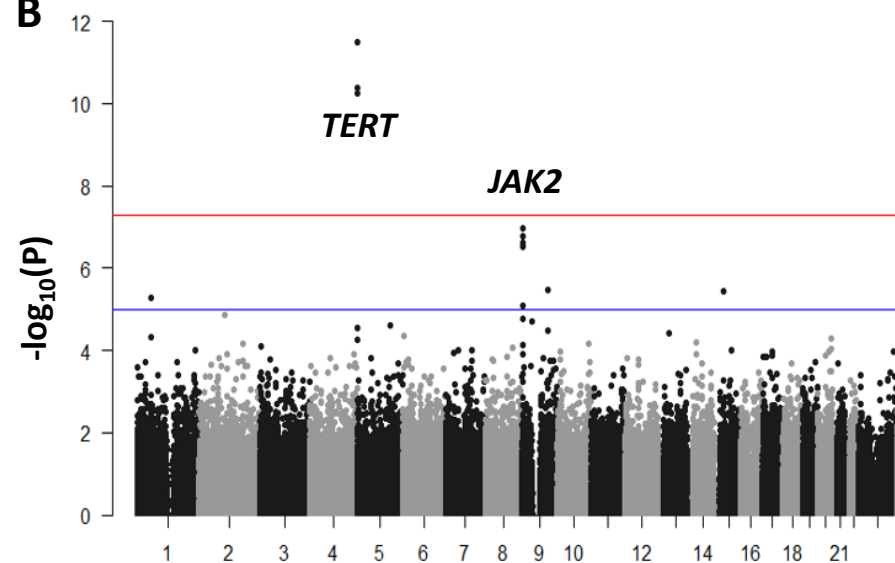
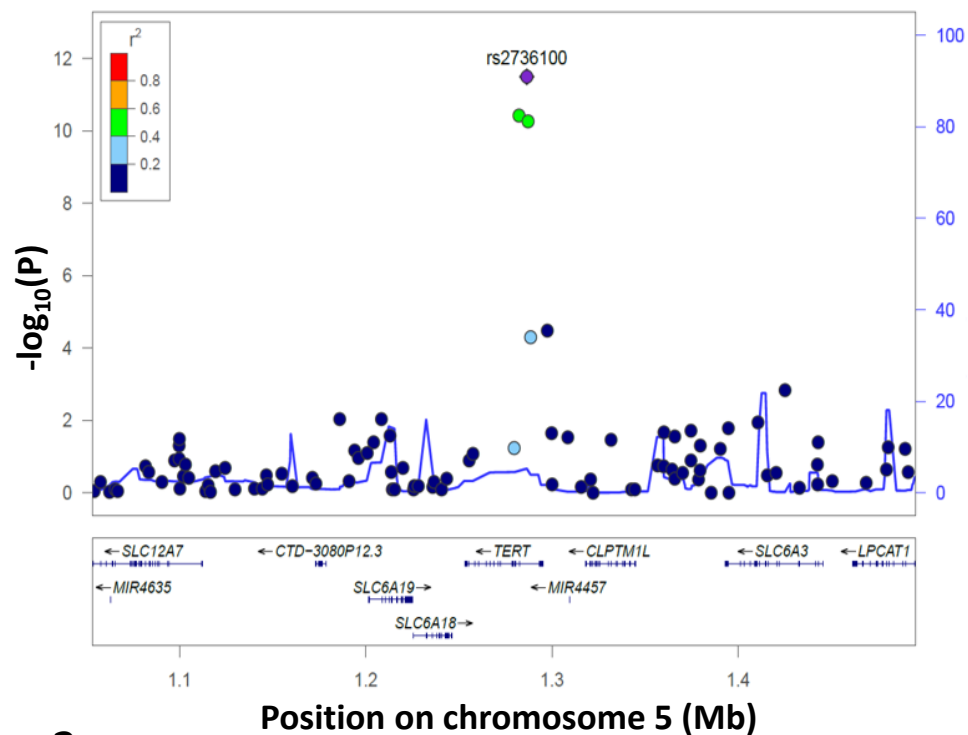
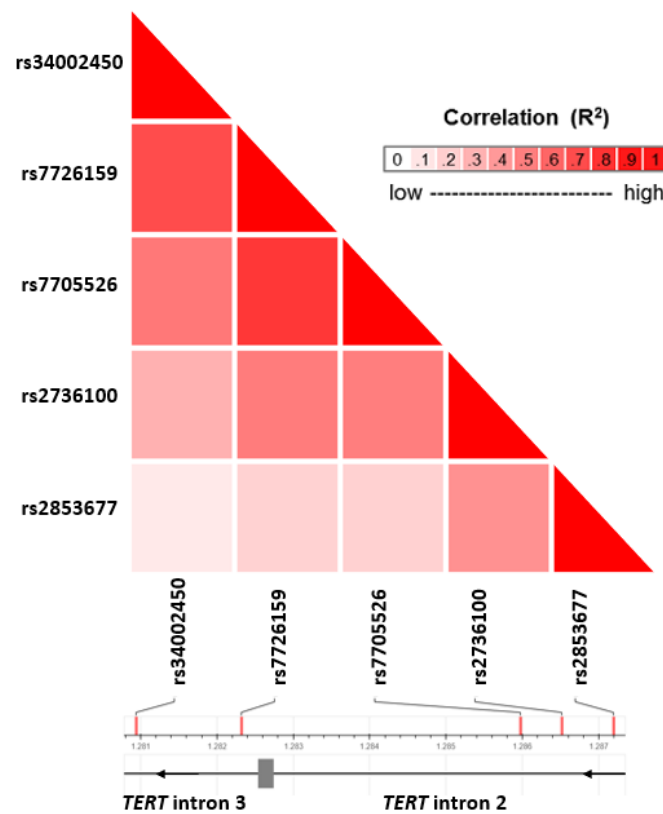


Figure 1

A**B****C****D****Figure 2**