

Genome-wide analyses of 200,453 individuals yields new insights into the causes and consequences of clonal hematopoiesis

Siddhartha P. Kar^{1,2*†}, Pedro M. Quiros^{3,4,5*†}, Muxin Gu^{3,4}, Tao Jiang⁶, Ryan Langdon^{1,2}, Vivek Iyer⁴, Clea Barcena^{3,4}, M.S. Vijayabaskar^{3,4}, Margarete A. Fabre^{3,4,7}, Paul Carter⁸, Stephen Burgess^{6,9}, and George S. Vassiliou^{3,4,7†}

1. MRC Integrative Epidemiology Unit, University of Bristol, Bristol, BS8 2BN, United Kingdom.
2. Section of Translational Epidemiology, Division of Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, BS8 2BN, United Kingdom.
3. Wellcome-MRC Cambridge Stem Cell Institute, University of Cambridge, Jeffrey Cheah Biomedical Centre, Cambridge, CB2 0AW, United Kingdom.
4. Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, United Kingdom.
5. Instituto de Investigación Sanitaria del Principado de Asturias, ISPA, 33011, Oviedo, Spain.
6. MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Strangeways Research Laboratory, Cambridge, CB1 8RN, United Kingdom.
7. Department of Haematology, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ, United Kingdom.
8. Division of Cardiovascular Medicine, Department of Medicine, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ, United Kingdom.
9. MRC Biostatistics Unit, University of Cambridge, Cambridge, CB2 0SR, United Kingdom.

* These authors contributed equally to this work.

† Correspondence: gsv20@cam.ac.uk ; siddhartha.kar@bristol.ac.uk ; pmquiros@ispasturias.es

Abstract

Clonal hematopoiesis (CH) is one of the most extensively studied somatic mutational phenomena, yet its causes and consequences remain poorly understood. We identify 10,924 individuals with CH amongst 200,453 whole-exome sequenced UK Biobank participants and use their linked genome-wide DNA genotypes to map the landscape of inherited predisposition to CH. We increase the number of European-ancestry genome-wide significant ($P < 5 \times 10^{-8}$) germline associations with CH from four to 14 and identify one new transcriptome-wide significant ($P < 3.2 \times 10^{-6}$) association. Genes at new loci implicate DNA damage repair (*PARP1*, *ATM*, and *CHEK2*), hematopoietic stem cell migration/homing (*CD164*), and myeloid oncogenesis (*SETBP1*) in CH development. Several associations were CH-subtype specific and, strikingly, variants at *TCL1A* and *CD164* had opposite associations with *DNMT3A*- versus *TET2*-mutant CH, mirroring recently reported differences in lifelong behavior of these two most common CH subtypes and proposing important roles for these loci in CH pathogenesis. Using Mendelian randomization, we show, amongst other findings, that smoking and longer leukocyte telomere length are causal risk factors for CH and demonstrate that genetic predisposition to CH increases risks of myeloproliferative neoplasia, several non-hematological malignancies, atrial fibrillation, and blood epigenetic age acceleration.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

52 Introduction

53

54 The pervasive effects of ageing and somatic mutation shape the landscape of human disease in later
55 life¹. A ubiquitous feature of ageing is the development of somatic mutation-driven clonal expansions
56 in aged tissues^{2,3}. In blood, somatic mutations that enhance cellular fitness of individual hematopoietic
57 stem cells (HSCs) and their progeny, give rise to the common age-related phenomenon of clonal
58 hematopoiesis (CH)⁴⁻⁷. CH becomes increasingly prevalent with advancing age⁴⁻⁶ and is associated
59 with an increased risk of hematological cancers^{4,5,8,9} and of some non-hematological conditions^{5,10,11}.
60 However, our understanding of the biological basis for these associations remains limited, as does our
61 ability to explain how CH driver mutations promote clonal expansion of mutant HSCs¹². In fact, whilst
62 CH is defined by its association with somatic mutations, its development is influenced by non-
63 mutation factors¹³⁻¹⁶ and by the heritable genome^{17,18}, in ways that remain poorly understood.

64

65 Insights into the causes and consequences of CH are confounded by its intimate relationship with
66 ageing. Moreover, even when robust associations are identified, their causality can be difficult to
67 establish. Here, we perform a comprehensive investigation of the genetic and phenotypic associations
68 of CH in 200,453 United Kingdom Biobank (UKB) participants, yielding a step change in our
69 understanding of CH pathogenesis. Our study reveals multiple new germline loci associated with CH,
70 including several that interact with specific CH subtypes, uncovers causal links between CH and diverse
71 pathological states across organ systems, and provides evidence for causal associations between
72 smoking and telomere length and CH risk, amongst a series of novel insights.

73

74 Results

75

76 *Prevalence of CH and its distribution by age and sex in the UKB*

77

78 To identify individuals with CH, we analyzed blood whole exome sequencing (WES) data from 200,453
79 UKB participants¹⁹ aged 38-72 years (Extended Data Fig. 1a-c). We called somatic mutations in 43 CH
80 genes (Supplementary Table 1) and filtered these against a predefined list of CH driver variants
81 (Supplementary Tables 2 and 3). This identified 11,697 mutations (Supplementary Table 4) in 10,924
82 individuals (UKB prevalence: 5.45%). *DNMT3A*, *TET2*, and *ASXL1* were most commonly involved (79%
83 of all mutations), followed by mutations in DNA damage response genes *PPM1D*, *TP53*, and *ATM*;
84 splicing factor genes *SRSF2* and *SF3B1*; *JAK2* and *GNB1* (Fig. 1a), in line with previous reports^{4,5,17}. Most
85 CH carriers (n=10,228) harbored one and some (n=696) 2-4 mutations; most of which were missense
86 variants dominated by cytosine-to-thymine (C>T) transitions (Extended Data Fig. 1d-f). The mean
87 variant allele fraction (VAF) was 0.12 and VAF distribution did not differ between mutation types
88 (Extended Data Fig. 1g and h). VAF distribution did differ between individual genes (Fig. 1a), although
89 some of this variation was probably influenced by variation in sequencing depth (Supplementary Table
90 1).

91

92 CH prevalence rose progressively with age ($P < 10^{-300}$; Fig. 1b), as did clone size measured by VAF
93 ($P = 8.2 \times 10^{-37}$; Extended Data Fig. 2a). Females and males were similarly affected with similar median
94 ages (Extended Data Fig. 2b). The age-related rise in prevalence differed between drivers: compared
95 to *DNMT3A*, mutations in *ATM* were observed earlier and those in *ASXL1*, *PPM1D*, *SF3B1*, and *SRSF2*
96 were observed later (Extended Data Fig. 2c). Furthermore, we noted significant differences in the
97 prevalence of different CH gene mutations between sexes, with *GNB1* and *DNMT3A* mutations more
98 frequent in females and *PPM1D*, *TP53*, *JAK2*, *SF3B1*, *ASXL1*, and *SRSF2* mutations more frequent in
99 males (Fig. 1c), reflecting their relative prevalence in myeloid malignancies²⁰.

100

101

102

103 *Associations between CH and traits/diseases prevalent at the time of blood sampling*

104

105 To identify associations between CH and traits or diseases prevalent at the time of enrolment to the
106 UKB, we performed regression analyses adjusting for age, sex, smoking status, WES batch and the first
107 ten genetic ancestry principal components. Individuals with CH showed higher average platelet,
108 leukocyte, reticulocyte, and neutrophil counts and red blood cell distribution width (RDW), but lower
109 eosinophil counts (Fig. 2a). These associations were more pronounced in individuals with large CH
110 clones (VAF \geq 0.1; Fig. 2a; Supplementary Table 5). *JAK2*-driven CH was associated with markedly higher
111 platelet counts, RDW and hemoglobin/hematocrit (HGB/HT) levels. In contrast, splicing factor-mutant
112 CH was associated with lower HGB/HT and higher mean red cell volume (MCV; Fig. 2a; Supplementary
113 Table 5). We also found that CH status was associated with lower levels of total and low-density
114 lipoprotein cholesterol (Fig. 2b; Supplementary Table 6). Advancing age increased the risk of CH by
115 6.7% per year (OR=1.07, 95%CI: 1.06-1.07, $P<10^{-300}$; Fig. 2c); and CH status was associated with
116 increased prevalence of hypertension, but not obesity or type 2 diabetes (T2D; Fig. 2c; Supplementary
117 Table 7). Also, individuals with CH were more likely to be current, past, or “ever” smokers, an
118 association that held true for different forms of CH and was strongest for *ASXL1*-mutant CH (Fig 2c
119 and 2d; Supplementary Table 7).

120

121 *Associations between CH and incident disease*

122

123 We next investigated relationships between CH at baseline and traits/diseases that developed
124 subsequently (Supplementary Table 8) and identified strong associations with incident myeloid
125 malignancies (MM) and associated sequelae (Extended Data Fig. 3a and Supplementary Table 9). The
126 association was strong for all MM subtypes and highest for chronic myelomonocytic leukemia (CMML;
127 Fig. 2e), whilst large clone CH increased the risk of MMs three-to-five-fold compared to small clone
128 CH (VAF $<$ 0.1; Extended Data Fig. 3b; Supplementary Table 10). *SF3B1* and *SRSF2* mutations conferred
129 very high risks of CMML and myelodysplastic syndromes (Extended Data Fig. 3c; Supplementary Table
130 10). CH was also associated with increased risks of Hodgkin’s and non-Hodgkin’s lymphomas and non-
131 hematological neoplasia, including lung, head and neck, kidney, bladder, colorectal, and stomach
132 cancers (Fig. 2f; Supplementary Table 11). The association of CH with lung adenocarcinoma was
133 consistently observed across large and small clones, and with *DNMT3A* and *ASXL1* mutations, whilst
134 the association with overall CH persisted in self-reported never-smokers (Extended Data Fig. 3d).

135

136 As CH was previously identified as a risk factor for ischemic cardiovascular disease (CVD)^{5,10,21}, we
137 examined the association in this much larger cohort (Fig. 2g; Supplementary Table 12). Using
138 multivariable models, we did not find a significant association between CH and ischemic CVD,
139 including coronary artery disease (CAD) and stroke; however, we did find significantly increased risks
140 of heart failure and a composite of all CVD conditions (Fig. 2g; Supplementary Table 12). Using a
141 bivariable model including age as the only other covariate, we also found a significant association with
142 atrial fibrillation, with an effect size estimate consistent with that in multivariable analysis (Extended
143 Data Fig. 3e; Supplementary Table 12). Our multivariable analyses also found significant associations
144 between CH and increased risk of death from any cause, malignant neoplasm, and hematological
145 neoplasm, whilst large clone CH was also associated with an increased risk of death due to the
146 composite of CVDs (Fig. 2h; Supplementary Table 13).

147

148 *Heritability and cell type-specific enrichment of polygenic susceptibility to CH*

149

150 To identify heritable determinants of CH risk, we performed a genome-wide association study (GWAS)
151 on the 184,121 individuals with genetically inferred European ancestry to identify common (minor
152 allele frequency (MAF) $>$ 1%) germline genetic variants predisposing to CH. In the GWAS, we compared
153 10,203 individuals with CH with 173,918 individuals without CH, after quality control of the germline

154 genotype data. Linkage disequilibrium score regression (LDSC)²² showed little evidence of inflation in
155 test statistics due to population structure, with an intercept of 1.009 and lambda genomic control
156 factor of 0.999. The narrow-sense (additive) heritability of CH was estimated at 3.57% (s.e.=0.85%).
157 We partitioned the heritability of CH across four major histone marks observed in 10 cell type groups
158 aggregated from 220 cell-type specific annotations²³ and identified strong enrichment of the polygenic
159 CH association signal in histone marks enriched in hematopoietic cells ($P=5.9 \times 10^{-5}$; Fig. 3a;
160 Supplementary Table 14). Next, we partitioned the heritability of CH across open chromatin state
161 regions in various hematopoietic progenitor cells and lineages^{23,24}. We found evidence of CH
162 heritability enrichment in accessible chromatin regions in HSCs, common lymphoid and myeloid
163 progenitors, multipotent and erythroid progenitors, and B cells (Fig. 3b; Supplementary Table 15).
164 Overall, these findings are in keeping with the intuitive assumption that the CH GWAS exerts its
165 greatest biological effect on HSC/progenitor populations.

166 167 *Germline genetic loci associated with overall CH susceptibility*

168
169 Linkage disequilibrium (LD)-based clumping of 9,715,652 common variants identified seven
170 independent ($r^2 < 0.05$) genome-wide significant loci (lead variant $P < 5 \times 10^{-8}$) associated with risk of
171 developing CH, including three previously reported¹⁷ European-ancestry CH loci: two at 5p15.33-*TERT*
172 and one at 3q25.33-*SMC4* (Fig. 4a; Supplementary Table 16). We identified a new top variant in the
173 5p15.33 region, rs2853677 ($P=2.4 \times 10^{-50}$), which was weakly correlated ($r^2=0.19$) with the previously
174 reported¹⁷ top variant, rs7705526 ($P=3.4 \times 10^{-44}$ in our analysis). Overall, there was evidence for three
175 independent ($r^2 < 0.05$) signals at 5p15.33 marked by lead variants rs2853677, rs13156167, and
176 rs2086132, the latter representing a new signal independent of the two previously published¹⁷ signals
177 rs7705526 and rs13167280. After approximate conditional analysis²⁵ (Supplementary Table 17)
178 conditioning on the three lead variants in the *TERT* region, the previously published top variant,
179 rs7705526, continued to remain genome-wide significant suggesting that it represented a fourth
180 signal in this region. Conditional analysis also highlighted the existence of a fifth independent
181 association at 5p15.33 marked by rs13356700 ~776 kb from *TERT* and ~34 kb from *EXOC3*
182 (Supplementary Table 17) that encodes an exocyst complex component implicated in arterial
183 thrombosis²⁶. The variant rs13356700 was in strong LD ($r^2=0.84$) with rs10072668 that is associated
184 with HGB/HT²⁷. At 3q25.33-*SMC4*, the previously reported¹⁷ top variant, rs1210060191, was not
185 captured in the UKB and our top association was rs12632224 ($P=2.3 \times 10^{-9}$). We also identified three
186 other novel genome-wide significant loci associated with overall CH susceptibility (Fig. 4a;
187 Supplementary Table 16): 4q35.1-*ENPP6* (rs13130545), 6q21-*CD164* (rs35452836), and 11q22.3-*ATM*
188 (rs11212666).

189 190 *CH GWAS stratified by gene and clone size, association heterogeneity, and rare variant associations*

191
192 Next, we investigated whether the development of certain CH subtypes may be affected by germline
193 variants. Thus, we performed GWAS for four additional CH traits – stratifying by the two main mutated
194 genes in CH, *DNMT3A* and *TET2*, and by clonal size, differentiating large and small clones. Focusing on
195 5,185 individuals with *DNMT3A* and 2,041 with *TET2* mutations and using the 173,918 individuals of
196 European ancestry without detectable CH as controls, we identified eight and three genome-wide
197 significant loci associated with *DNMT3A*- and *TET2*-mutant CH, respectively (Figs. 4b and 4c;
198 Supplementary Tables 18 and 19). We replicated the only previously published European-ancestry CH
199 risk locus associated with *DNMT3A*-CH at 14q32.13-*TCL1A*. The overall CH loci at 5p15.33-*TERT* (signals
200 with lead variants rs2853677, rs13156167, and rs7705526), 3q25.33-*SMC4*, 6q21-*CD164*, and
201 11q22.3-*ATM* were also genome-wide significant for *DNMT3A*-mutant CH. We also found two novel
202 loci for *DNMT3A*-CH marked by lead variants rs138994074 at 1q42.12-*PARP1* and rs8088824 at
203 18q12.3-*SETBP1* (Fig. 4b; Supplementary Table 18). The three *TET2*-CH associated loci included the
204 lead variant rs2736100 at 5p15.33-*TERT*, that was moderately correlated ($r^2=0.44$) with the overall CH

205 lead variant rs2853677 in the same region. The other two risk loci, both new in the context of *TET2*-
206 CH, were at lead variants rs10131341 (14q32.13-*TCL1A*) and rs79633204 (7q32.2-*TMEM209*; Fig. 4c;
207 Supplementary Table 19). Notably, the A allele of rs10131341 had opposite associations with *TET2*-CH
208 (OR=1.28, $P=6.8 \times 10^{-10}$) versus *DNMT3A*-CH (OR=0.87, $P=6.4 \times 10^{-8}$). A trend for opposite effects at
209 14q32.13-*TCL1A* was also observed in a previous study¹⁷, but did not achieve genome-wide
210 significance for *TET2*-CH.

211
212 When comparing 4,049 individuals with large or 6,154 individuals with small clones against 173,918
213 controls of European ancestry without CH, we found that the overall CH loci at 5p15.33-*TERT* and
214 3q25.33-*SMC4* were associated at genome-wide significance with large clone CH (Fig. 4d;
215 Supplementary Table 20), while 5p15.33-*TERT* and 6q21-*CD164* were associated with small clone CH.
216 For small clone CH risk, we also identified a previously unreported locus marked by rs72755524 at
217 5p13.3 in a region with several lincRNAs (Fig. 4e; Supplementary Table 21). Additional signals
218 suggested by approximate conditional analysis at each locus identified in this study are listed in
219 Supplementary Table 17. Examining heterogeneity of associations across the five CH traits using forest
220 plots (Extended Data Fig. 4) revealed that in addition to 14q32.13-*TCL1A*, the lead alleles at 6q21-
221 *CD164* also had opposite effects on *DNMT3A*- versus *TET2*-CH. In addition, the lead variants at 6q21-
222 *CD164* and 5p13.3-*LINC02064* were associated with small, but not large, clones while the association
223 at 7q32.2-*TMEM209* was highly specific to *TET2*-CH.

224
225 Finally, in addition to our common variant GWAS, we performed a more focused scan to explore rare
226 variant (MAF: 0.2%-1%) associations in each of three CH traits that included >5,000 European-ancestry
227 individuals with CH (i.e., overall CH, *DNMT3A*-CH, and small clone CH; each compared to 173,918
228 controls) declaring associations significant at $P < 10^{-9}$. This identified one new locus at 22q12.1-*CHEK2*
229 where the T allele (frequency=0.3%) of lead variant rs62237617 was perfectly correlated ($r^2=1$) with
230 the 1100delC *CHEK2* protein-truncating allele (rs555607708) and conferred a large increase in risk of
231 *DNMT3A* mutation-associated CH (OR=4.1, 95%CI: 2.7-6.1, $P=6.3 \times 10^{-12}$). The *CHEK2* c.1100delC
232 frameshift mutation or its tagging variant rs62237617 are known to be associated with
233 myeloproliferative neoplasms (MPNs) and *JAK2* V617F-driven CH (though not genome-wide significant
234 for either trait)¹⁸, elevated white blood cell counts and plateletcrit²⁷, as well as risk of prostate and
235 breast cancers^{28,29}. The *DNMT3A*-CH risk increasing alleles in the *CHEK2* and *PARP1* regions were also
236 associated with later age at menopause in a recent analysis³⁰, suggesting a role for inhibition of DNA
237 damage sensing and apoptosis in both CH and reproductive ageing³¹.

238
239 *Genetic relationship between hematological chromosomal mosaicism and CH due to gene mutation*

240
241 It is not known whether the germline genetic architecture underlying predisposition to CH due to
242 individual gene mutations is similar to that underlying the risk of CH due to mosaic chromosomal
243 alterations (mCAs). We used data from a recent GWAS of blood mCAs³² to answer this question and
244 found that 13 of 19 unique lead variants identified for the five gene-mutant CH traits (overall,
245 *DNMT3A*, *TET2*-, and large and small clone CH) were associated with hematological mCA risk at $P < 10^{-4}$
246 (Supplementary Table 22). Notably, for our lead variants rs2296312 (14q32.13-*TCL1A*) and rs8088824
247 (18q12.3-*SETBP1*), the alleles conferring increased *DNMT3A*-CH risk reduced hematological mCA risk
248 (Supplementary Table 22). At the genome-wide level we found a correlation between overall CH and
249 mCAs ($r_g=0.44$, s.e.=0.21, $P=0.037$) using LDSC²². Further, a phenome-wide scan^{33,34} showed that
250 several newly identified lead variants in our analyses were associated with multiple blood cell counts
251 and traits (Supplementary Table 23).

252
253
254
255

256 *Gene-level associations and network analyses*

257

258 We supplemented our GWAS with gene-level association tests for each of our five CH traits using two
259 complementary methods: multi-marker analysis of genomic annotation (MAGMA) and a
260 transcriptome-wide association study (TWAS) using blood-based *cis* gene expression quantitative trait
261 locus data on 31,684 individuals³⁵ and summary-based Mendelian randomization (SMR) coupled with
262 the heterogeneity in dependent instruments colocalization test³⁶. Both approaches converged on a
263 new locus at 6p21.1, associated at gene-level genome-wide significance ($P_{\text{MAGMA}} < 2.6 \times 10^{-6}$,
264 $P_{\text{SMR}} < 3.2 \times 10^{-6}$) with *DNMT3A*-mutant CH and marked by *CRIP3* ($P_{\text{MAGMA}} = 3.4 \times 10^{-7}$, $P_{\text{SMR}} = 6.6 \times 10^{-7}$; Fig. 3a;
265 Supplementary Tables 24 and 25). While *CRIP3* is the only 6p21.1 gene to reach gene-level genome-
266 wide significance in both MAGMA and SMR, we did find sub-threshold evidence for association
267 between *SRF* or *ZNF318* in the same region and *DNMT3A*-mutant CH (Fig. 5a). Of note, *SRF* encodes
268 the serum response factor that is known to regulate HSC adhesion³⁷ while *ZNF318* is an occasional
269 somatic driver gene for CH³⁸. More globally, protein-protein interaction (PPI) network analysis³⁹ using
270 proteins encoded by the 57 genes with $P_{\text{MAGMA}} < 0.001$ in the overall CH analysis (Supplementary Table
271 24) as “seeds”, identified the largest sub-network (Fig. 5b) as encompassing 13/57 proteins with major
272 hub nodes highlighted as TERT, PARP1, ATM, and SMC4. This was consistent with the emerging theme
273 that key genes at sub-threshold GWAS loci for the same trait are often part of interconnected
274 biological networks^{40,41}. The sub-threshold genes identified by MAGMA that encoded protein hubs in
275 this network included *FANCF* (DNA repair pathway) and *PTCH1* (hedgehog signaling; Fig. 5b), both
276 implicated in the pathogenesis of acute myeloid leukemia^{42,43} and *GNAS*, a somatic driver of CH⁴⁴. The
277 CH sub-network (seeds and non-seed interacting proteins) was significantly enriched for several
278 pathways of relevance to common disease including DNA repair, cell cycle regulation, telomere
279 maintenance, and platelet homeostasis (Supplementary Table 26).

280

281 *Functional target gene prioritization at CH risk loci*

282

283 In order to prioritize putative functional target genes at the $P_{\text{lead-variant}} < 5 \times 10^{-8}$ loci identified by our
284 GWAS of five CH traits, we combined gene-level genome-wide significant results from MAGMA and
285 SMR (Supplementary Tables 24 and 25) with five other lines of evidence: PPI network hub status of
286 the gene (Supplementary Table 27), variant-to-gene searches of the Open Targets database⁴⁵ for lead
287 variants, overlap between fine-mapped variants^{46,47} (Supplementary Table 28) and (i) gene bodies, (ii)
288 regions with accessible chromatin correlated with nearby gene expression in hematopoietic
289 progenitor cells^{24,48-50}, and (iii) missense variant annotations^{51,52} (Supplementary Table 29). Genes
290 nominated by at least two of these approaches are listed in Fig. 5c. The genes nominated by the largest
291 number of approaches, and representing the most likely targets, were *SMC4*, *ENPP6*, *TERT*, *CD164*,
292 *ATM*, *PARP1*, *TCL1A*, *SETBP1*, and *TMEM209*.

293

294 Among the newly identified loci, *CD164* codes for Sialomucin core protein 24, a cell adhesion molecule
295 that regulates HSC adhesion, proliferation, and migration^{53,54}. Lead variant rs138994074 at 1q42.12
296 was strongly correlated ($r^2 = 0.93$) with rs1136410, a missense germline mutation in *PARP1*
297 (Supplementary Table 29) wherein the G allele, which is protective for *DNMT3A*-CH, leads to a
298 missense variant (p.Val762Ala) in the catalytic domain of its protein product associated with reduced
299 Poly (ADP-ribose) polymerase-1 activity⁵⁵. While *SETBP1* was only nominated by one approach (Open
300 Targets⁴⁵) and was the only gene nominated at 18q12.3, its nomination is strengthened by the fact
301 that somatic *SETBP1* mutations are recognized drivers of myeloid malignancies^{56,57}.

302

303 *Mendelian randomization (MR) to uncover the causes and consequences of CH*

304

305 We integrated several large GWAS datasets (Supplementary Tables 30 and 31) and used two-sample
306 inverse-variance-weighted MR⁵⁸ to appraise putative causes and consequences of CH. Genetically-

307 predicted smoking initiation⁵⁹ was associated with overall CH risk (OR=1.15, 95%CI: 1.05-1.25,
308 $P=2.2 \times 10^{-3}$). Point estimates of the effect size were consistent in direction across MR analyses for
309 *DNMT3A*, *TET2*, and large and small clone CH (Fig. 6a; Supplementary Table 32), with the largest odds
310 ratio observed for large clone CH (OR=1.24). We also appraised the roles of leukocyte telomere length
311 (LTL)⁶⁰, alcohol use⁵⁹, adiposity⁶¹, genetic liability to T2D⁶², major circulating lipids⁶³, blood-based
312 epigenetic aging phenotypes⁶⁴, blood cell counts and indices²⁷, and circulating cytokines and growth
313 factors⁶⁵ as potential risk factors for CH using MR (Fig. 6; Supplementary Tables 32, 33, and 34 for full
314 results, including sensitivity analyses). Genetically predicted longer LTL was associated with increased
315 overall CH risk (OR=1.56, 95%CI: 1.25-1.93, $P=5.7 \times 10^{-5}$), an association that was also seen with
316 *DNMT3A*-, *TET2*-, and large and small clone CH (Fig. 6b; Supplementary Table 32). We found that
317 higher genetically predicted BMI was associated with increased risk of large clone CH (OR=1.15, 95%CI:
318 1.01-1.31, $P=0.029$). Genetically elevated circulating apolipoprotein B levels were associated with
319 increased (OR=1.18, 95%CI: 1.01-1.36, $P=0.032$; Fig. 6c), whilst genetically predicted alcohol use was
320 associated with decreased (OR=0.46, 95%CI: 0.25-0.83, $P=0.010$) risk of *TET2*-CH. Among cytokines,
321 genetically-elevated circulating macrophage inflammatory protein 1a, a regulator of myeloid
322 differentiation and HSC numbers⁶⁶, was associated with risk of *DNMT3A*-CH (OR=1.13, 95%CI 1.03-
323 1.23, $P=7.1 \times 10^{-3}$; Supplementary Table 34).

324
325 We used independent ($r^2 < 0.001$) variants associated with overall, *DNMT3A*, *TET2*, and large and small
326 clone CH at $P < 10^{-5}$ as genetic instruments for each of these traits and assessed their associations with
327 outcomes (Supplementary Tables 31, 35, and 36 for full results, including sensitivity analyses). Since
328 more variants were available at $P < 5 \times 10^{-8}$ for overall and for *DNMT3A* CH, we also examined the
329 consistency of associations when using genome-wide (GWS; $P < 5 \times 10^{-8}$) and sub-genome-wide
330 significant (sub-GWS; $P < 10^{-5}$) instruments for these two traits. Using the sub-GWS instrument, genetic
331 liability to overall CH had the largest associations (Fig. 7a) with MPN risk⁴⁸ (OR=1.99, 95%CI: 1.23-3.23,
332 $P=5.4 \times 10^{-3}$), intrinsic epigenetic age acceleration⁶⁴ (IEAA, which represents a core characteristic of
333 HSCs⁶⁷; beta= 0.39, 95%CI: 0.08-0.69, $P=0.01$) and the blood-based Hannum epigenetic clock⁶⁴ (beta=
334 0.27, 95%CI: 0.04-0.49, $P=0.02$). Larger associations were observed when using the GWS instrument
335 (Fig. 7a) and the direction of these was consistent when evaluating genetic liability to *DNMT3A*, *TET2*,
336 and large and small clone CH as exposures (Supplementary Tables 35 and 36). Genetic liability to CH
337 conferred increased risks of lung⁶⁸, prostate⁶⁹, ovarian⁷⁰, oral cavity/pharyngeal⁷¹, and endometrial
338 cancers⁷² (Fig. 7; Supplementary Table 35) with the strongest associations observed between overall
339 CH and lung (OR=1.17, 95%CI: 1.05-1.29, $P=2.9 \times 10^{-3}$); *DNMT3A*-CH and prostate (OR=1.08, 95%CI:
340 1.03-1.13, $P=8.6 \times 10^{-4}$), ovarian (OR=1.07, 95%CI: 1.01-1.12, $P=0.015$), and oral cavity/pharyngeal
341 (OR=1.24, 95%CI: 1.07-1.44, $P=4.4 \times 10^{-3}$); and *TET2*-CH and endometrial (OR=1.05, 95%CI: 1.00-1.09,
342 $P=0.033$) cancers. MR analyses did not support causal risk-conferring associations between genetic
343 liability to CH and CAD⁷³, ischemic stroke⁷⁴, and heart failure⁷⁵ with similar lack of evidence across
344 gene-specific and clone size-specific CH, and GWS instrument analyses (Fig. 7; Supplementary Table
345 35). However, we did uncover an association between genetic liability to overall CH or *DNMT3A*-CH
346 and atrial fibrillation⁷⁶ risk (OR=1.09, 95%CI: 1.04-1.15, $P=4.9 \times 10^{-4}$ for overall CH with the GWS
347 instrument; Supplementary Table 35). Among cytokines and growth factors⁶⁵, genetic liability to
348 overall CH was associated with elevated circulating stem cell growth factor beta (beta= 0.19; 95%CI:
349 0.07-0.30, $P=1.1 \times 10^{-3}$). MR analyses also revealed bidirectional associations between CH phenotypes
350 and several blood cell counts and traits²⁷, suggesting a shared underlying genetic liability to CH and
351 pan-blood cell proliferation (Figs. 6b and 7; Supplementary Tables 33 and 35). Finally, we found little
352 evidence to support an association between genetic liability to CH and LTL (Supplementary Table 36),
353 indicating that longer LTL was a cause rather than a consequence of CH.

354
355
356
357

358 Discussion

359

360 We present a large observational and genetic epidemiological analysis of CH and report a series of
361 novel insights into the causes and consequences of this common aging-associated phenomenon. We
362 increase the number of germline associations with CH in European-ancestry populations from four¹⁷
363 to 14, reveal heterogeneity of associations by CH driver gene and clone size, and implicate putative
364 new CH susceptibility genes, including *CD164*, *ATM* and *SETBP1*, through functional annotation. We
365 also demonstrate that the CH GWAS signal is enriched at epigenetic marks specific to the
366 hematopoietic system. The robustness of our GWAS analysis is further affirmed by our replication of
367 previous European ancestry-specific CH associations¹⁷, the consistency of our estimates of CH
368 heritability with previous reports^{17,77}, and the fact that many of our lead variants are associated with
369 related traits^{27,32,60,78}.

370

371 New CH risk loci included the *PARP1* coding variant rs1136410, where the G allele is protective for
372 *DNMT3A*-CH and associated with reduced catalytic activity⁵⁵ suggesting that this most common form
373 of CH may be vulnerable to PARP inhibition, in keeping with the observed synergy between PARP and
374 DNMT inhibitors⁷⁹. At 14q32.13-*TCL1A*, we replicate the reported association with *DNMT3A*-CH¹⁷ and
375 identify a new genome-wide significant association with *TET2*-CH. Strikingly, however, we found that
376 the association operates in the opposite direction for *TET2*-CH, versus *DNMT3A*-CH. This inverse
377 relationship is tantalizing in light of recent observations that ageing has different effects on the
378 dynamics of these two forms of CH, resulting in *TET2* CH becoming more prevalent than *DNMT3A* CH
379 in those aged over 80 years^{80,81}. Also notable in this light, is the finding of an association at the *CD164*
380 locus with *DNMT3A*, and a trend in the opposite direction for *TET2*-CH. As *CD164* is expressed in the
381 earliest HSCs⁵³ and encodes an important regulator of HSC adhesion^{54,82}, this proposes that HSC
382 migration and homing may play important roles in CH pathogenesis. The reciprocal relationship of
383 both *TCL1A* and *CD164* with the two main CH subtypes, suggests that their expression needs to be
384 tightly regulated to prevent the development of one or other subtype of CH, making these loci
385 important targets for hijack by the effects of somatic mutations.

386

387 The rich phenotypic data captured by the UKB, coupled with our genetic analysis of CH and external
388 GWAS datasets, enabled us to explore associations of CH using multivariable regression and
389 interrogate, at scale, potential causal relationships between CH and its putative risk factors and
390 consequences using MR. This highlighted for the first time that smoking and longer telomere length
391 are causal risk factors for CH. These associations were valid across multiple CH subtypes and, in the
392 case of smoking, corroborated by observational estimates. We also reveal that not only is genetic
393 predisposition to CH causally associated with MPN risk, but it also increases the risk of lung, prostate,
394 ovarian, oral/pharyngeal, and endometrial cancers. In these analyses, the use of two-sample MR
395 protected against potential reverse causality arising from cancer therapy-induced selection pressure
396 on hematopoietic clones⁸³. These MR results suggest that genetic liability to CH may be a biomarker
397 for development of cancer elsewhere in the body. An analogous relationship has previously been
398 identified by MR for the association of genetic predisposition to Y chromosome loss in blood and solid
399 tumor risk³¹.

400

401 We investigated the recently identified association of CH with blood-based epigenetic clocks⁸⁴, using
402 bi-directional MR and show that this association is likely to be causal in the direction from CH to
403 epigenetic age acceleration. We also showed that genetic predisposition to CH was associated with
404 elevated circulating levels of stem cell growth factor beta, a secreted sulfated glycoprotein that
405 regulates primitive hematopoietic progenitor cells⁸⁵. Finally, we unraveled a previously unreported
406 association between genetic liability to CH and atrial fibrillation risk, which was also supported by our
407 observational analysis. However, unlike previous reports based on significantly smaller sample
408 numbers^{5,10,21}, we did not find evidence in observational and MR analyses to support an association

409 between CH and CAD or ischemic stroke risk. However, our MR analyses indicated that higher BMI
410 and circulating apolipoprotein B levels were associated with *TET2* and large clone CH risks,
411 respectively, with apolipoprotein B being the key causal lipid risk factor for CAD^{63,86}. These associations
412 taken together with the fact that age and smoking are strong risk factors for CH raise the possibility
413 that previously reported associations of CH with CAD and stroke risks may suffer from residual
414 confounding.

415
416 Collectively, our findings substantially illuminate the landscape of inherited susceptibility to CH and
417 provide new insights into the causes and consequences of CH with implications for human health and
418 ageing.

419 420 **Methods**

421 422 *Study population and exome sequence data*

423
424 The United Kingdom Biobank (UKB) is a prospective longitudinal study containing in-depth genetic and
425 health information from half a million UK participants. For this study, we have selected 200,453
426 individuals (200k) who had whole exome sequencing (WES) data available (age range: 38-72, median
427 age: 58; 55% female; 83% White British). WES was generated in two batches, the first of approximately
428 50,000 samples (50k)⁸⁷ and the second comprising an additional 150,000 samples (150k)¹⁹. Exomes
429 were captured using the IDT xGen Exome Research Panel v1.0 including supplemental probes; a
430 different IDT v1.0 oligo lot was used for each batch. Multiplexed samples were sequenced with dual-
431 indexed 75x75 bp paired-end reads on the Illumina NovaSeq 6000 platform using S2 (50k samples)
432 and S4 (150k samples) flow cells. The 50k samples were firstly computed using FE protocol and
433 reprocessed later to match the second batch of 150k sequences that were processed using a new
434 improved unified OQFE pipeline. As the initial 50k samples were sequenced on S2 flow cells and with
435 a different IDT v1.0 oligo lot than the remaining 150k samples, which were sequenced on S4 flow cells,
436 we included the WES batch as a covariate in downstream analyses.

437
438 The UK Biobank study has been approved by the North West Multicentre Research Ethics Committee
439 (11/NW/0382). All participants provided written informed consent. The current study has been
440 conducted under approved UK Biobank application numbers 56844 and 29202.

441 442 *Whole exome sequence data processing, CH mutation calling and filtering*

443
444 CRAM files generated by the OQFE pipeline were obtained from UKB (Fields 23143-23144;
445 www.ukbiobank.ac.uk). Variant-calling on WES data from 200,453 individuals was performed using
446 Mutect2, Genome Analysis Toolkit (GATK) version 4.1.8.1⁸⁸. Briefly, Mutect2 was run in “tumor-only”
447 mode with default parameters, over the exon intervals of 43 genes previously associated with CH
448 (Supplementary Table 1). To filter out potential germline variants we used a population reference of
449 germline variants generated from 1000 Genomes Project (1000GP)⁸⁹ and the Genome Aggregation
450 Database (gnomAD)⁹⁰. All resources were obtained from the GATK Best practices repository ([gs://gatk-
best-practices/somatic-hg38](https://gatk-best-practices.github.io/somatic-hg38/)). Raw variants called by Mutect2 were filtered out with *FilterMutectCalls*
451 using the estimated prior probability of a reading orientation artefact generated by
452 *LearnReadOrientationModel* (GATK v.4.1.8.1). Putative variants flagged as ‘PASS’ using
453 *FilterMutectCalls* or flagged as ‘germline’ if present at least 2 times with the ‘PASS’ flag in other
454 samples were selected for filtering. Gene annotation was performed using Ensembl Variant Effect
455 Predictor (VEP) (v.102)⁹¹. We required variants with a minimum number of alternate reads of 2,
456 evidence of the variant on both forward and reverse strand, a minimum depth of 7 reads for SNVs and
457 10 reads for short indels and substitutions and a minor allele frequency (MAF) lower than 0.001
458 (according to 1000GP phase 3 and gnomAD r2.1). For new variants, not previously described in the
459

460 Catalogue of Somatic Mutations in Cancer (COSMIC; v.91)⁹² nor in the Database of Single Nucleotide
461 Polymorphisms (dbSNP; build 153)⁹³, we used a minimum allele count per variant of 4, and a MAF
462 lower than 5×10^{-5} . From resulting variants, we selected those that: i) are included in a list of recurring
463 hotspots mutations associated with CH and myeloid cancer (Supplementary Table 2); ii) have been
464 reported as somatic mutations in hematological cancers at least 7 times in COSMIC; or iii) met the
465 inclusion criteria of a predefined list of putative CH variants, previously described^{17,77} (Supplementary
466 Table 3). We included previous variants flagged as germline by *FilterMutectCalls* if: 1) the number of
467 cases in the cohort flagged as germline were lower than the ones flagged as PASS; and 2) at least one
468 of the cases had a $P < 0.001$ for a one-sided exact binomial test, where the null hypothesis was that the
469 number of alternative reads supporting the mutation were 50% of the total number of reads (95% for
470 copy number equal to one), except for hotspot mutations that were all included. For the final list, we
471 excluded all variants not present in COSMIC nor in the list of hotspots that had a MAF equal or higher
472 than 5×10^{-5} and either the mean variant allele fraction (VAF) of all cases was higher than 0.2 or the
473 maximum VAF was lower than 0.1. Frameshift, nonsense, and splice-site mutations not present in
474 COSMIC nor in the hotspot list were further excluded if for each variant none of the cases had a
475 $P < 0.001$ for a one-sided exact binomial test. A complete list of filtered variants is provided in
476 Supplementary Table 4.

477

478 *Trait selection and modelling for the conventional observational multivariable regression analyses*

479

480 Phenotypes were downloaded in December 2020 and individual traits were pulled out from the whole
481 phenotype file. Cancer, metabolic and cardiovascular disease (CVD) traits were generated combining
482 individual traits and diagnosis dates based on disease definitions (Supplementary Table 8). For each
483 definition of disease, the first diagnosis event that occurred in each trait was selected. Baseline was
484 defined as the date of sample collection. The prevalent cases are those identified before the baseline,
485 while incidence was defined as the events that occurred after the baseline. Unless specified, all
486 regression models included age, sex, smoking status, WES batch and the first ten ancestry principal
487 components as covariates. Blood cell counts and biochemical traits were \log_{10} transformed and
488 analyzed using a linear regression model, including the assessment center as covariate and, in the case
489 of cholesterol and cholesterol species, the use of cholesterol lowering medication. Individuals with
490 myeloid malignancies or hematological neoplasms at baseline were excluded from the analysis. For
491 cancer, CVD and death risk, we performed a time-to-event regression analysis using the Cox
492 proportional hazards model. The cancer/CVD/death event was used as an outcome and CH was
493 considered as the exposure in these analyses. For CVD and death risk analyses, we also included body
494 mass index, high-density lipoprotein cholesterol, low-density lipoprotein cholesterol, triglycerides,
495 type 2 diabetes status, and hypertension status as covariates. Individuals with myeloid or other
496 malignant neoplasms at baseline were excluded from all previous analyses. For associations between
497 International Statistical Classification of Diseases and Related Health Problems 10th Revision (ICD-10)
498 codes and CH status, a logistic regression model was used including age, sex, WES batch and the first
499 ten ancestry principal components as covariates. Analyses were performed over the selected ICD-10
500 codes corresponding to diseases conditions (A to N), symptoms, signs, and abnormal clinical and
501 laboratory findings (R) and factors influencing health status (Z). All analyses were performed using glm
502 (R stats package v.4.0.2) and coxph (R survival package v.3.2-11) functions.

503

504 *Germline genotype data processing and genome-wide association analyses*

505

506 Germline genotype data used were from the UKB release that contained the full set of variants
507 imputed into the Haplotype Reference Consortium⁹⁴ and 1000GP⁸⁹ reference panels and genotyped
508 on the UK BiLEVE Axiom Array or UKB Axiom Array⁹⁵. Derivation of the analytic sample for UK Biobank
509 of individuals of European ancestries followed quality control (QC) steps described previously²⁷: after
510 filtering genetic variants (call rate $\geq 99\%$, imputation quality info score > 0.9 , Hardy-Weinberg

511 equilibrium P -value $\geq 10^{-5}$) and participants (removal of genetic sex mismatches), we excluded
512 participants having non-European ancestries (self-report or inferred by genetics) or excess
513 heterozygosity (>3 standard deviations from the mean), and included only one of each set of related
514 participants (third-degree relatives or closer). After QC, we were left with 10,203 individuals with CH
515 and 173,918 individuals without CH. The subset with CH included 5,185 and 2,041 individuals with
516 *DNMT3A* and *TET2*-mutant CH, respectively, and 4,049 and 6,154 individuals with large ($VAF \geq 0.1$) and
517 small ($VAF < 0.1$) clone size CH, respectively. Association analyses were performed using non-
518 infinitesimal linear mixed models implemented in BOLT-LMM⁹⁶ with age at baseline, sex, and first 10
519 genetic principal components included as covariates.

520
521 Statistically independent lead variants for each CH phenotype were defined using linkage
522 disequilibrium (LD)-based clumping with an r^2 threshold of 0.05 applied across all genotyped and
523 imputed variants with $P < 5 \times 10^{-8}$, imputation quality score > 0.6 , and $MAF > 1\%$. This was implemented
524 using the FUMA pipeline⁹⁷. For the rare variant association scan, we used more stringent cut-offs of
525 $P < 10^{-9}$ and imputation quality score > 0.8 to define lead variants but did not require LD-clumping since
526 only one such association was identified. Approximate conditional analysis conditioning on the
527 common ($MAF > 1\%$) lead variants was performed using the `--cojo-cond` flag in the Genome-wide
528 Complex Trait Analysis (GCTA) v1.93 tool^{25,98}.

529
530 *Linkage disequilibrium score regression (LDSC)*

531
532 We used LDSC²² to estimate the narrow-sense heritability of CH on the liability scale assuming the
533 population prevalence of CH to be 10% (based on the prevalence of CH in the UKB “200k” cohort as
534 shown in Fig. 1b) and constraining the LDSC intercept to 1. The intercept, which in its unconstrained
535 form protects from bias due to population stratification, was constrained to 1 to provide more precise
536 estimates given that there was little evidence of inflation in test statistics due to population structure
537 in unconstrained analysis (unconstrained intercept estimated as 1.009 (s.e.=0.0067) and lambda
538 genomic control factor of 0.999). We used the pre-computed 1000 Genomes phase 3 European
539 ancestry reference panel LD score data set downloaded from alkesgroup.broadinstitute.org/LDSCORE
540 for heritability estimation. We used the same LD scores and the `--rg` flag in LDSC to estimate the
541 genetic correlation between the CH and mosaic chromosomal alteration GWAS summary statistics³².
542 Cell-type group partitioned heritability analysis was performed as described in
543 github.com/bulik/ldsc/wiki/Partitioned-Heritability using LD scores partitioned across 220 cell-type-
544 specific annotations that were divided into 10 groups as previously described²³: central nervous
545 system, cardiovascular, kidney, adrenal/pancreas, gastrointestinal, connective/bone,
546 immune/hematopoietic, skeletal muscle, liver, and other. Each of the 10 groups contained cell-type-
547 specific annotations for four histone marks: H3K9ac, H3K27ac, H3K4me1, and H3K4me3²³. We also
548 used LD scores annotated as previously described⁹⁹ based on open chromatin state (Assay for
549 Transposase-Accessible Chromatin (ATAC)-seq) profiling by Corces et al.²⁴ in various hematopoietic
550 progenitor cells and lineages at different stages of differentiation.

551
552 *Gene-based and transcriptome-wide association studies, and network analyses*

553
554 We undertook genome-wide gene-level association analyses using two complementary approaches.
555 First, we used multi-marker analysis of genomic annotation (MAGMA) that involves mapping germline
556 variants to the genes they overlap, accounting for LD between variants, and performing a statistical
557 multi-marker association test¹⁰⁰. Second, we performed a transcriptome-wide association study
558 (TWAS) using blood-based *cis* gene expression quantitative trait locus (eQTL) data on 31,684
559 individuals³⁵ and summary-based Mendelian randomization (SMR) coupled with the heterogeneity in
560 dependent instruments (HEIDI) colocalization test to identify germline genetic associations with CH
561 risk mediated via the transcriptome³⁶. The gene-level genome-wide significance threshold in the

562 MAGMA analyses was set at $P=2.6 \times 10^{-6}$ to account for testing 19,064 genes and for SMR was set at
563 $P=3.2 \times 10^{-6}$ after adjustment for testing 15,672 genes. Further, only genes with SMR $P < 3.2 \times 10^{-6}$ and
564 HEIDI $P > 0.05$ were declared genome-wide significant in the SMR analyses since the HEIDI $P > 0.05$
565 strongly suggests colocalization of the GWAS and eQTL signals for a given gene³⁶. The NetworkAnalyst
566 3.0³⁹ webtool available at www.networkanalyst.ca was used for network analysis. All genes with $P < 10^{-3}$
567 in each MAGMA analysis for overall, *DNMT3A* and *TET2*-mutant, and large and small clone CH were
568 used as input. The protein-protein interactome selected was STRING v10¹⁰¹ with the recommended
569 parameters (confidence score cut-off of 900 and requirement for experimental evidence to support
570 the protein-protein interaction). The largest possible network was constructed from the seed
571 genes/proteins and the interactome proteins as previously described³⁹. Hub nodes were defined as
572 nodes with degree centrality ≥ 10 (i.e., nodes with at least 10 edges or connections to other proteins in
573 the network as a measure of its importance in the network and consequently its biology). Pathway
574 analysis of this largest network was conducted using the enrichment tool built into the
575 NetworkAnalyst webtool and with the Reactome pathway repository¹⁰².

576 577 *Fine-mapping and target gene prioritization*

578
579 We fine-mapped the lead variant signals identified by the FUMA LD-clumping pipeline using the
580 Probabilistic Identification of Causal Single Nucleotide Polymorphisms (PICS2) algorithm^{46,47} to identify
581 candidate causal variants most likely to underpin each association. The PICS2 algorithm and webtool
582 (pics2.ucsf.edu) computes the likelihood that each variant in LD with the lead variant is the true causal
583 variant in the region by leveraging the fact that for variants associated merely due to LD, the strength
584 of association scales asymptotically with correlation to the true causal variant⁴⁶. We only retained
585 variants with a PICS2 probability of 1% or more in our final list of fine-mapped candidate causal
586 variants. We overlapped these fine-mapped variants with gene body annotations as previously
587 described⁴⁸ using GENCODE release 33¹⁰³ (build 37) annotations after removing ribosomal protein
588 genes (code and data adapted from github.com/sankaranlab/mpn-gwas). Fine-mapped variants were
589 also overlapped with ATAC-seq peaks across 16 hematopoietic progenitor cell populations and ATAC-
590 RNA count correlations calculated using Pearson coefficients for hematopoietic progenitor cell RNA
591 counts of genes within 1 Mb of the ATAC peaks were used to identify putative target genes of fine-
592 mapped variants that overlapped ATAC-seq peaks. This pipeline has been used and described
593 extensively before^{24,48-50}, and we adapted the code and data for the pipeline from
594 github.com/sankaranlab/mpn-gwas. We also looked up the SIFT⁵¹ and PolyPhen⁵² scores for these
595 fine-mapped variants using the SNPnexus v4 web-based annotation tool (www.snp-nexus.org/v4)¹⁰⁴
596 to identify coding variants with predicted functional consequences. Finally, we used the Open Targets
597 Genetics resource⁴⁵ (genetics.opentargets.org) to identify the most likely target gene of the lead
598 variant at each locus as per Open Targets and used this in our omnibus target gene prioritization
599 scheme described below.

600
601 In order to prioritize putative target genes at the $P_{\text{lead-variant}} < 5 \times 10^{-8}$ loci identified by our GWAS of
602 overall CH, *DNMT3A*-CH, *TET2*-CH and large/small clone size CH, we combined gene-level genome-
603 wide significant results from (1) MAGMA and (2) SMR with (3) protein-protein interaction network
604 hub status of the gene, (4) variant-to-gene searches of the Open Targets database for lead variants,
605 and overlap between fine-mapped variants and (5) gene bodies, (6) regions with accessible chromatin
606 (ATAC-seq peaks) across 16 hematopoietic progenitor cell populations that were also correlated with
607 nearby gene expression (RNA-seq) in the same cell populations, and (7) missense variant annotations
608 from SIFT and PolyPhen. Genes nominated by at least two of the seven approaches were listed (except
609 where only one of the seven methods nominated a single gene in a region in which case that gene was
610 listed) and the genes nominated by the largest number of approaches represented the most likely
611 targets at each locus.
612

613 *Phenome-wide association scan for lead variants*

614

615 We used PhenoScanner V2^{33,34} available at www.phenoscanter.medschl.cam.ac.uk with catalogue set
616 to “diseases & traits”, p-value set to “5E-8”, proxies set to “EUR” and r^2 set to “0.8” to search for
617 published phenome-wide associations between our lead variants or variants in strong linkage
618 disequilibrium ($r^2 > 0.8$) with the lead variants and other diseases and traits.

619

620 *Mendelian randomization analysis*

621

622 Mendelian randomization (MR)^{105,106} uses germline variants as instrumental variables to proxy an
623 exposure or potential risk factor and evaluate evidence for a causal effect of the exposure or potential
624 risk factor on an outcome. Due to the random segregation and independent assortment of alleles at
625 meiosis, MR estimates are less susceptible to bias from confounding factors as compared to
626 conventional observational epidemiological studies. As the germline genome cannot be influenced by
627 the environment after conception or by preclinical disease, MR estimates are also less susceptible to
628 bias due to reverse causation. MR estimates represent the association between genetically predicted
629 levels of exposures or risk factors and outcomes, as compared to conventional observational
630 epidemiological estimates, which represent direct associations of the exposure or risk factor levels
631 with outcomes. Effect allele harmonization across GWAS summary statistics datasets followed by two-
632 sample Mendelian randomization analyses were performed using the TwoSampleMR v0.5.6 R
633 package⁵⁸. The CH phenotypes were considered as both exposures (to identify consequences of
634 genetic liability to CH) and outcomes (to identify risk factors for CH). When considering CH phenotypes
635 as outcomes, germline variants associated with putative risk factors or exposures at $P < 5 \times 10^{-8}$ were
636 used as genetic instruments for the risk factors/exposures, except for the appraisal of circulating
637 cytokines and growth factors⁶⁵ wherein variants associated with cytokines/growth factors at $P < 10^{-5}$
638 were used as instruments. Inverse-variance weighted analysis¹⁰⁷ was the primary analytic approach
639 with pleiotropy-robust sensitivity analyses carried out using the MR-Egger¹⁰⁸ and weighted median¹⁰⁹
640 methods. A full list of external GWAS data sources used for MR analyses is provided in Supplementary
641 Tables 30 and 31.

642

643 **Acknowledgements**

644

645 This work was funded by a joint grant from the Leukemia and Lymphoma Society (RTF6006-19) and
646 the Rising Tide Foundation for Clinical Cancer Research (CCR-18-500), and by the Wellcome Trust
647 (WT098051). SPK is supported by a United Kingdom Research and Innovation (UKRI) Future Leaders
648 Fellowship (MR/T043202/1) and leads the somatic genomics theme of the Integrative Cancer
649 Epidemiology Programme at the University of Bristol that is funded by Cancer Research UK
650 (C18281/A29019). PMQ is funded by the Miguel Servet Program (CP20/00130). MAF is funded by a
651 Wellcome Clinical Research Fellowship (WT098051). RL is supported by Cancer Research UK
652 (C18281/A29019). PC is supported by a British Heart Foundation Clinical Training Research Fellowship.
653 SB is supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal
654 Society (204623/Z/16/Z). GSV is supported by a Cancer Research UK Senior Cancer Fellowship
655 (C22324/A23015) and work in his lab is also funded by the European Research Council, Kay Kendall
656 Leukaemia Fund, Blood Cancer UK, and the Wellcome Trust. This research was conducted using the
657 UK Biobank resource under applications 56844 and 29202. We thank the participants and
658 investigators involved in the UK Biobank resource and in the other genome-wide association studies
659 cited in this work who collectively made this research possible.

660

661

662

663

664 **Ethics Declarations**

665

666 *Competing Interests*

667

668 GSV is a consultant to STRM.BIO and AstraZeneca.

669

670 **Author Contributions**

671

672 SPK, PMQ, and GSV conceived, designed, and supervised the study. SPK and PMQ carried out data
673 analyses and generated tables and figures. MG, MSV, and MAF helped with mutation calling and
674 filtering. TJ and SB performed genome-wide association analyses. RL assisted with Mendelian
675 randomization analyses. VI helped with UK Biobank data access and handling. SB and PC advised on
676 Mendelian randomization analyses. CB and PC helped with UK Biobank trait selection and filtering.
677 SPK, PMQ, and GSV drafted the manuscript with inputs from all authors. All authors approved the final
678 version of the paper.

679

680 **Data Availability**

681

682 Individual-level UK Biobank data can be requested via application to the UK Biobank
683 (<https://www.ukbiobank.ac.uk>). The CH call set will be returned to the UK Biobank to enable
684 individual-level data linkage for approved UK Biobank applications.

685

686 **Code Availability**

687

688 Code used in this study is available at https://github.com/pmquiros/CH_UKBiobank and
689 <https://github.com/siddhartha-kar/clonal-hematopoiesis>.

690

691 **Figures, Extended Data Figures, and Figure Legends can be found from pages 21 to 31**

692

693 **References**

694

695 1. Zhang, L. & Vijg, J. Somatic Mutagenesis in Mammals and Its Implications for Human Disease and Aging. *Annu. Rev.*

696 *Genet.* **52**, 397–419 (2018).

697 2. Martincorena, I. & Campbell, P. J. Somatic mutation in cancer and normal cells. *Science* **349**, 1483–1489 (2015).

698 3. Kakiuchi, N. & Ogawa, S. Clonal expansion in non-cancer tissues. *Nat. Rev. Cancer* **21**, 239–256 (2021).

699 4. Genovese, G. *et al.* Clonal hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *N. Engl. J. Med.*

700 **371**, 2477–2487 (2014).

701 5. Jaiswal, S. *et al.* Age-related clonal hematopoiesis associated with adverse outcomes. *N. Engl. J. Med.* **371**, 2488–2498
702 (2014).

703 6. McKerrell, T. *et al.* Leukemia-associated somatic mutations drive distinct patterns of age-related clonal hemopoiesis.

704 *Cell Rep.* **10**, 1239–1245 (2015).

705 7. Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion and malignancies. *Nat. Med.* **20**,

706 1472–1478 (2014).

707 8. Abelson, S. *et al.* Prediction of acute myeloid leukaemia risk in healthy individuals. *Nature* **559**, 400–404 (2018).

- 708 9. Desai, P. *et al.* Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nat. Med.* **24**, 1015–1023
709 (2018).
- 710 10. Jaiswal, S. *et al.* Clonal Hematopoiesis and Risk of Atherosclerotic Cardiovascular Disease. *N. Engl. J. Med.* **377**, 111–121
711 (2017).
- 712 11. Dorsheimer, L. *et al.* Association of Mutations Contributing to Clonal Hematopoiesis With Prognosis in Chronic Ischemic
713 Heart Failure. *JAMA Cardiol.* **4**, 25–33 (2019).
- 714 12. Jaiswal, S. & Ebert, B. L. Clonal hematopoiesis in human aging and disease. *Science* **366**, eaan4673 (2019).
- 715 13. Coombs, C. C. *et al.* Therapy-Related Clonal Hematopoiesis in Patients with Non-hematologic Cancers Is Common and
716 Associated with Adverse Clinical Outcomes. *Cell Stem Cell* **21**, 374-382.e4 (2017).
- 717 14. Gibson, C. J. *et al.* Clonal Hematopoiesis Associated With Adverse Outcomes After Autologous Stem-Cell
718 Transplantation for Lymphoma. *J. Clin. Oncol.* **35**, 1598–1605 (2017).
- 719 15. Meisel, M. *et al.* Microbial signals drive pre-leukaemic myeloproliferation in a Tet2-deficient host. *Nature* **557**, 580–584
720 (2018).
- 721 16. Yoshizato, T. *et al.* Somatic Mutations and Clonal Hematopoiesis in Aplastic Anemia. *N. Engl. J. Med.* **373**, 35–47 (2015).
- 722 17. Bick, A. G. *et al.* Inherited causes of clonal haematopoiesis in 97,691 whole genomes. *Nature* **586**, 763–768 (2020).
- 723 18. Hinds, D. A. *et al.* Germ line variants predispose to both JAK2 V617F clonal hematopoiesis and myeloproliferative
724 neoplasms. *Blood* **128**, 1121–1128 (2016).
- 725 19. Szustakowski, J. D. *et al.* Advancing human genetics research and drug discovery through exome sequencing of the UK
726 Biobank. *Nat. Genet.* **53**, 942–948 (2021).
- 727 20. De-Morgan, A., Meggendorfer, M., Haferlach, C. & Shlush, L. Male predominance in AML is associated with specific
728 preleukemic mutations. *Leukemia* **35**, 867–870 (2021).
- 729 21. Bick, A. G. *et al.* Genetic Interleukin 6 Signaling Deficiency Attenuates Cardiovascular Risk in Clonal Hematopoiesis.
730 *Circulation* **141**, 124–131 (2020).
- 731 22. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association
732 studies. *Nat. Genet.* **47**, 291–295 (2015).
- 733 23. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary
734 statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 735 24. Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia
736 evolution. *Nat. Genet.* **48**, 1193–1203 (2016).
- 737 25. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants
738 influencing complex traits. *Nat. Genet.* **44**, 369–375, S1-3 (2012).

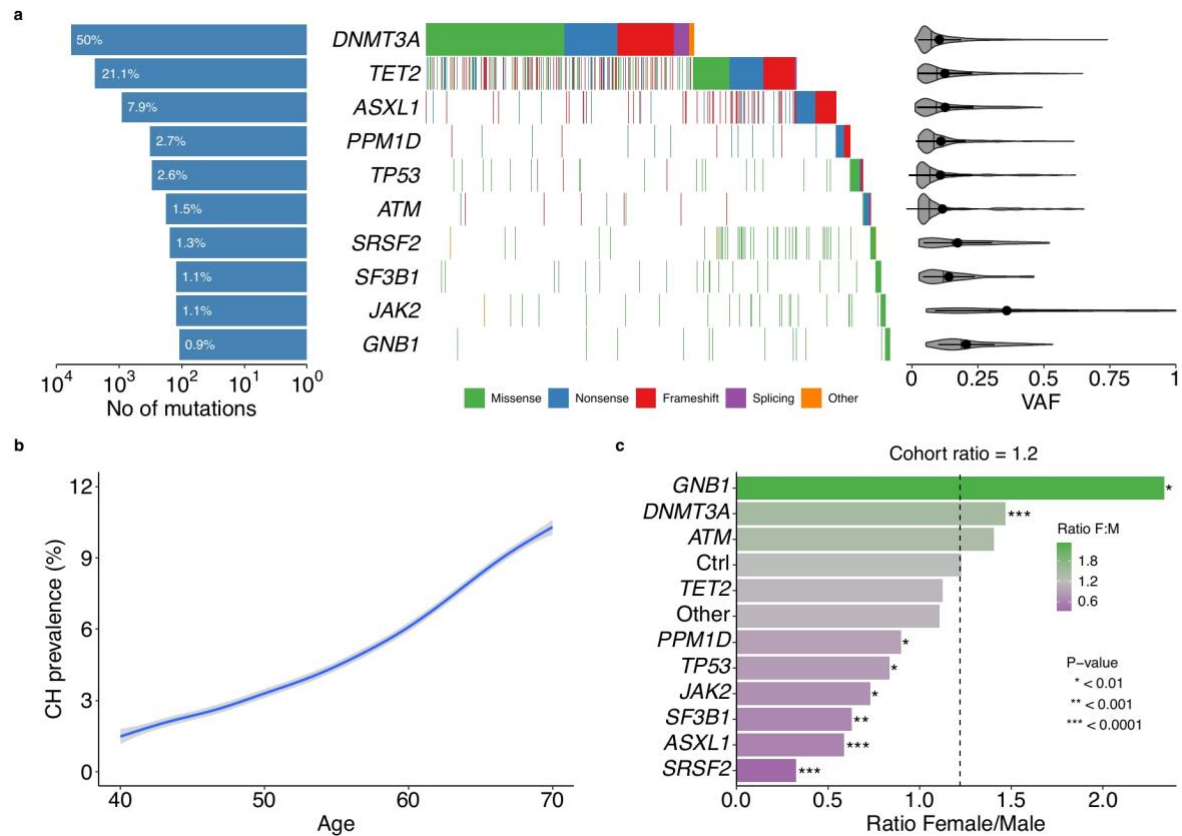
- 739 26. Walsh, T. G. *et al.* Loss of the exocyst complex component EXOC3 promotes hemostasis and accelerates arterial
740 thrombosis. *Blood Adv.* **5**, 674–686 (2021).
- 741 27. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell*
742 **167**, 1415–1429.e19 (2016).
- 743 28. Conti, D. V. *et al.* Trans-ancestry genome-wide association meta-analysis of prostate cancer identifies new susceptibility
744 loci and informs genetic risk prediction. *Nat. Genet.* **53**, 65–75 (2021).
- 745 29. Schmidt, M. K. *et al.* Age- and Tumor Subtype-Specific Breast Cancer Risk Estimates for CHEK2*1100delC Carriers. *J.*
746 *Clin. Oncol.* **34**, 2750–2760 (2016).
- 747 30. Ruth, K. S. *et al.* Genetic insights into biological mechanisms governing human ovarian ageing. *Nature* **596**, 393–397
748 (2021).
- 749 31. Thompson, D. J. *et al.* Genetic predisposition to mosaic Y chromosome loss in blood. *Nature* **575**, 652–657 (2019).
- 750 32. Zekavat, S. M. *et al.* Hematopoietic mosaic chromosomal alterations increase the risk for diverse types of infection. *Nat.*
751 *Med.* **27**, 1012–1024 (2021).
- 752 33. Kamat, M. A. *et al.* PhenoScanner V2: an expanded tool for searching human genotype-phenotype associations.
753 *Bioinformatics* **35**, 4851–4853 (2019).
- 754 34. Staley, J. R. *et al.* PhenoScanner: a database of human genotype-phenotype associations. *Bioinformatics* **32**, 3207–3209
755 (2016).
- 756 35. Vösa, U. *et al.* *Unraveling the polygenic architecture of complex traits using blood eQTL metaanalysis.* (2018)
757 doi:10.1101/447367.
- 758 36. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.*
759 **48**, 481–487 (2016).
- 760 37. Ragu, C. *et al.* The transcription factor Srf regulates hematopoietic stem cell adhesion. *Blood* **116**, 4464–4473 (2010).
- 761 38. Zink, F. *et al.* Clonal hematopoiesis, with and without candidate driver mutations, is common in the elderly. *Blood* **130**,
762 742–752 (2017).
- 763 39. Zhou, G. *et al.* NetworkAnalyst 3.0: a visual analytics platform for comprehensive gene expression profiling and meta-
764 analysis. *Nucleic Acids Res.* **47**, W234–W241 (2019).
- 765 40. Cowen, L., Ideker, T., Raphael, B. J. & Sharan, R. Network propagation: a universal amplifier of genetic associations. *Nat.*
766 *Rev. Genet.* **18**, 551–562 (2017).
- 767 41. Fernández-Tajes, J. *et al.* Developing a network view of type 2 diabetes risk pathways through integration of genetic,
768 genomic and functional data. *Genome Med.* **11**, 19 (2019).
- 769 42. Tischkowitz, M. *et al.* Bi-allelic silencing of the Fanconi anaemia gene FANCF in acute myeloid leukaemia. *Br. J.*
770 *Haematol.* **123**, 469–471 (2003).

- 771 43. Lim, Y. *et al.* Integration of Hedgehog and mutant FLT3 signaling in myeloid leukemia. *Sci. Transl. Med.* **7**, 291ra96
772 (2015).
- 773 44. Ostrander, E. L. *et al.* The GNASR201C mutation associated with clonal hematopoiesis supports transplantable
774 hematopoietic stem cell activity. *Exp. Hematol.* **57**, 14–20 (2018).
- 775 45. Ghossaini, M. *et al.* Open Targets Genetics: systematic identification of trait-associated genes using large-scale
776 genetics and functional genomics. *Nucleic Acids Res.* **49**, D1311–D1320 (2021).
- 777 46. Taylor, K. E., Ansel, K. M., Marson, A., Criswell, L. A. & Farh, K. K.-H. PICS2: Next-generation fine mapping via
778 probabilistic identification of causal SNPs. *Bioinformatics* btab122 (2021).
- 779 47. Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* **518**, 337–343
780 (2015).
- 781 48. Bao, E. L. *et al.* Inherited myeloproliferative neoplasm risk affects haematopoietic stem cells. *Nature* **586**, 769–775
782 (2020).
- 783 49. Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human
784 Hematopoietic Differentiation. *Cell* **173**, 1535–1548.e16 (2018).
- 785 50. Ulirsch, J. C. *et al.* Interrogation of human hematopoiesis at single-cell and single-variant resolution. *Nat. Genet.* **51**,
786 683–693 (2019).
- 787 51. Vaser, R., Adusumalli, S., Leng, S. N., Sikic, M. & Ng, P. C. SIFT missense predictions for genomes. *Nat. Protoc.* **11**, 1–9
788 (2016).
- 789 52. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248–249
790 (2010).
- 791 53. Pellin, D. *et al.* A comprehensive single cell transcriptional landscape of human hematopoietic progenitors. *Nat.*
792 *Commun.* **10**, 2395 (2019).
- 793 54. Zannettino, A. C. *et al.* The sialomucin CD164 (MGC-24v) is an adhesive glycoprotein expressed by human
794 hematopoietic progenitors and bone marrow stromal cells that serves as a potent negative regulator of hematopoiesis.
795 *Blood* **92**, 2613–2628 (1998).
- 796 55. Wang, X.-G., Wang, Z.-Q., Tong, W.-M. & Shen, Y. PARP1 Val762Ala polymorphism reduces enzymatic activity. *Biochem.*
797 *Biophys. Res. Commun.* **354**, 122–126 (2007).
- 798 56. Makishima, H. *et al.* Somatic SETBP1 mutations in myeloid malignancies. *Nat. Genet.* **45**, 942–946 (2013).
- 799 57. Piazza, R. *et al.* Recurrent SETBP1 mutations in atypical chronic myeloid leukemia. *Nat. Genet.* **45**, 18–24 (2013).
- 800 58. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human phenome. *eLife* **7**,
801 (2018).

- 802 59. Liu, M. *et al.* Association studies of up to 1.2 million individuals yield new insights into the genetic etiology of tobacco
803 and alcohol use. *Nat. Genet.* **51**, 237–244 (2019).
- 804 60. Codd, V. *et al.* *Polygenic basis and biomedical consequences of telomere length variation*. (2021)
805 doi:10.1101/2021.03.23.21253516.
- 806 61. Pulit, S. L. *et al.* Meta-analysis of genome-wide association studies for body fat distribution in 694 649 individuals of
807 European ancestry. *Hum. Mol. Genet.* **28**, 166–174 (2019).
- 808 62. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-
809 specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
- 810 63. Richardson, T. G. *et al.* Evaluating the relationship between circulating lipoprotein lipids and apolipoproteins with risk
811 of coronary heart disease: A multivariable Mendelian randomisation analysis. *PLoS Med.* **17**, e1003062 (2020).
- 812 64. McCartney, D. L. *et al.* Genome-wide association studies identify 137 genetic loci for DNA methylation biomarkers of
813 aging. *Genome Biol.* **22**, 194 (2021).
- 814 65. Ahola-Olli, A. V. *et al.* Genome-wide Association Study Identifies 27 Loci Influencing Concentrations of Circulating
815 Cytokines and Growth Factors. *Am. J. Hum. Genet.* **100**, 40–50 (2017).
- 816 66. Staversky, R. J. *et al.* The Chemokine CCL3 Regulates Myeloid Differentiation and Hematopoietic Stem Cell Numbers.
817 *Sci. Rep.* **8**, 14691 (2018).
- 818 67. Horvath, S. & Raj, K. DNA methylation-based biomarkers and the epigenetic clock theory of ageing. *Nat. Rev. Genet.* **19**,
819 371–384 (2018).
- 820 68. Wang, Y. *et al.* Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. *Nat. Genet.* **46**, 736–741
821 (2014).
- 822 69. Schumacher, F. R. *et al.* Association analyses of more than 140,000 men identify 63 new prostate cancer susceptibility
823 loci. *Nat. Genet.* **50**, 928–936 (2018).
- 824 70. Phelan, C. M. *et al.* Identification of 12 new susceptibility loci for different histotypes of epithelial ovarian cancer. *Nat.*
825 *Genet.* **49**, 680–691 (2017).
- 826 71. Lesueur, C. *et al.* Genome-wide association analyses identify new susceptibility loci for oral cavity and pharyngeal
827 cancer. *Nat. Genet.* **48**, 1544–1550 (2016).
- 828 72. O’Mara, T. A. *et al.* Identification of nine new susceptibility loci for endometrial cancer. *Nat. Commun.* **9**, 3166 (2018).
- 829 73. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-analysis of coronary artery
830 disease. *Nat. Genet.* **47**, 1121–1130 (2015).
- 831 74. Malik, R. *et al.* Multiancestry genome-wide association study of 520,000 subjects identifies 32 loci associated with
832 stroke and stroke subtypes. *Nat. Genet.* **50**, 524–537 (2018).

- 833 75. Shah, S. *et al.* Genome-wide association and Mendelian randomisation analysis provide insights into the pathogenesis
834 of heart failure. *Nat. Commun.* **11**, 163 (2020).
- 835 76. Nielsen, J. B. *et al.* Biobank-driven genomic discovery yields new insight into atrial fibrillation biology. *Nat. Genet.* **50**,
836 1234–1239 (2018).
- 837 77. Fabre, M. A. *et al.* Concordance for clonal hematopoiesis is limited in elderly twins. *Blood* **135**, 269–273 (2020).
- 838 78. Vuckovic, D. *et al.* The Polygenic and Monogenic Basis of Blood Traits and Diseases. *Cell* **182**, 1214–1231.e11 (2020).
- 839 79. Muvarak, N. E. *et al.* Enhancing the Cytotoxic Effects of PARP Inhibitors with DNA Demethylating Agents - A Potential
840 Therapy for Cancer. *Cancer Cell* **30**, 637–650 (2016).
- 841 80. Rossi, M. *et al.* Clinical relevance of clonal hematopoiesis in the oldest-old population. *Blood* blood.2021011320 (2021).
- 842 81. Fabre, M. A. *et al.* *The longitudinal dynamics and natural history of clonal haematopoiesis.* (2021)
843 doi:10.1101/2021.08.12.455048.
- 844 82. Forde, S. *et al.* Endolyn (CD164) modulates the CXCL12-mediated migration of umbilical cord blood CD133+ cells. *Blood*
845 **109**, 1825–1833 (2007).
- 846 83. Bolton, K. L. *et al.* Cancer therapy shapes the fitness landscape of clonal hematopoiesis. *Nat. Genet.* **52**, 1219–1226
847 (2020).
- 848 84. Nachun, D. *et al.* Clonal hematopoiesis associated with epigenetic aging and clinical outcomes. *Aging Cell* **20**, e13366
849 (2021).
- 850 85. Hiraoka, A. *et al.* Stem cell growth factor: in situ hybridization analysis on the gene expression, molecular
851 characterization and in vitro proliferative activity of a recombinant preparation on primitive hematopoietic progenitor
852 cells. *Hematol. J.* **2**, 307–315 (2001).
- 853 86. Zuber, V. *et al.* High-throughput multivariable Mendelian randomization analysis prioritizes apolipoprotein B as key
854 lipid risk factor for coronary artery disease. *Int. J. Epidemiol.* **50**, 893–901 (2021).
- 855 87. Van Hout, C. V. *et al.* Exome sequencing and characterization of 49,960 individuals in the UK Biobank. *Nature* **586**, 749–
856 756 (2020).
- 857 88. Auwera, G. A. V. de & O'Connor, B. D. *Genomics in the cloud: using Docker, GATK, and WDL in Terra.* (O'Reilly, 2020).
- 858 89. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 859 90. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**,
860 434–443 (2020).
- 861 91. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
- 862 92. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res.* **47**, D941–D947 (2019).
- 863 93. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308–311 (2001).
- 864 94. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).

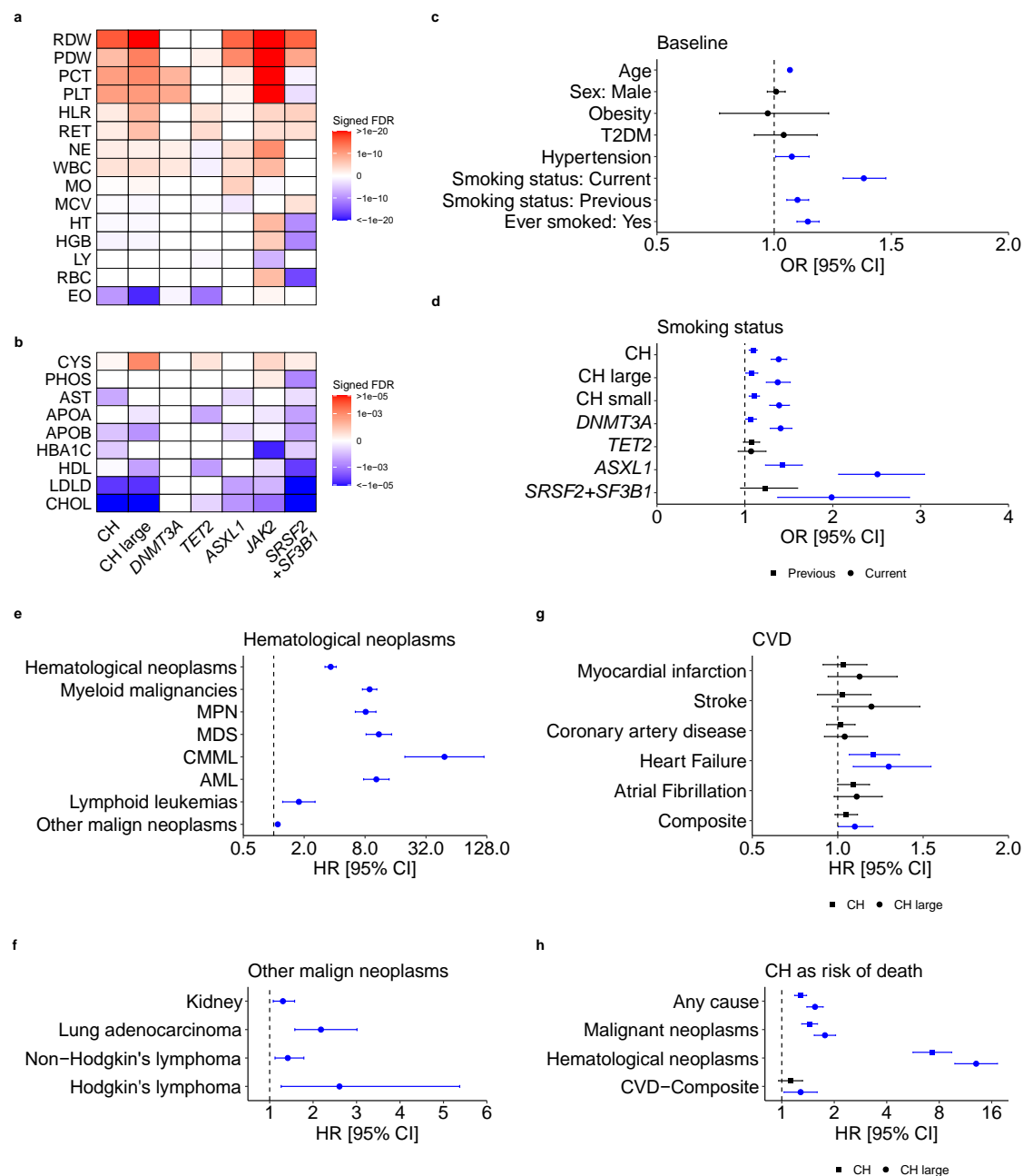
- 865 95. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
- 866 96. Loh, P.-R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* **47**,
867 284–290 (2015).
- 868 97. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and annotation of genetic
869 associations with FUMA. *Nat. Commun.* **8**, 1826 (2017).
- 870 98. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum.*
871 *Genet.* **88**, 76–82 (2011).
- 872 99. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell
873 types. *Nat. Genet.* **50**, 621–629 (2018).
- 874 100. de Leeuw, C. A., Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: generalized gene-set analysis of GWAS data.
875 *PLoS Comput. Biol.* **11**, e1004219 (2015).
- 876 101. Szklarczyk, D. *et al.* STRING v10: protein-protein interaction networks, integrated over the tree of life. *Nucleic*
877 *Acids Res.* **43**, D447–452 (2015).
- 878 102. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–D503 (2020).
- 879 103. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**,
880 D766–D773 (2019).
- 881 104. Oscanoa, J. *et al.* SNPnexus: a web server for functional annotation of human genome sequence variation (2020
882 update). *Nucleic Acids Res.* **48**, W185–W192 (2020).
- 883 105. Davey Smith, G. & Hemani, G. Mendelian randomization: genetic anchors for causal inference in epidemiological
884 studies. *Hum. Mol. Genet.* **23**, R89–98 (2014).
- 885 106. Davies, N. M., Holmes, M. V. & Davey Smith, G. Reading Mendelian randomisation studies: a guide, glossary, and
886 checklist for clinicians. *BMJ* **362**, k601 (2018).
- 887 107. Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis with multiple genetic variants
888 using summarized data. *Genet. Epidemiol.* **37**, 658–665 (2013).
- 889 108. Bowden, J., Davey Smith, G. & Burgess, S. Mendelian randomization with invalid instruments: effect estimation
890 and bias detection through Egger regression. *Int. J. Epidemiol.* **44**, 512–525 (2015).
- 891 109. Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent Estimation in Mendelian Randomization with
892 Some Invalid Instruments Using a Weighted Median Estimator. *Genet. Epidemiol.* **40**, 304–314 (2016).
- 893
894



895
896
897
898
899
900
901
902
903
904
905

Fig. 1: Characterization of CH in the UK Biobank.

a, Composite plot summarizing mutations in the ten most common driver genes in 10,924 individuals with CH. Each column in the waterfall plot represents a single individual, with mutation types color-coded. Bars on the left quantify mutations per gene as a percentage of all CH mutations identified. Violin plots on the right show the distribution of variant allele fractions (VAFs), with vertical lines represent the median and dots with horizontal lines the mean \pm standard deviation. **b**, Prevalence of CH in the cohort with advancing age. The blue line represents the smoothed model fitted to a generalized additive model with 95% confidence interval (CI; grey shadow). **c**, Bar plot showing the female to males (F:M) ratio of CH carriers with mutations in the ten most common driver genes. "Other" represents the remaining driver genes grouped together and "Ctrl" the ratio for individuals without CH. Dotted vertical line shows the F:M ratio observed in the full cohort (F:M=1.2). *P*-values are from a Chi-square test comparing the distribution for each gene to "Ctrl".

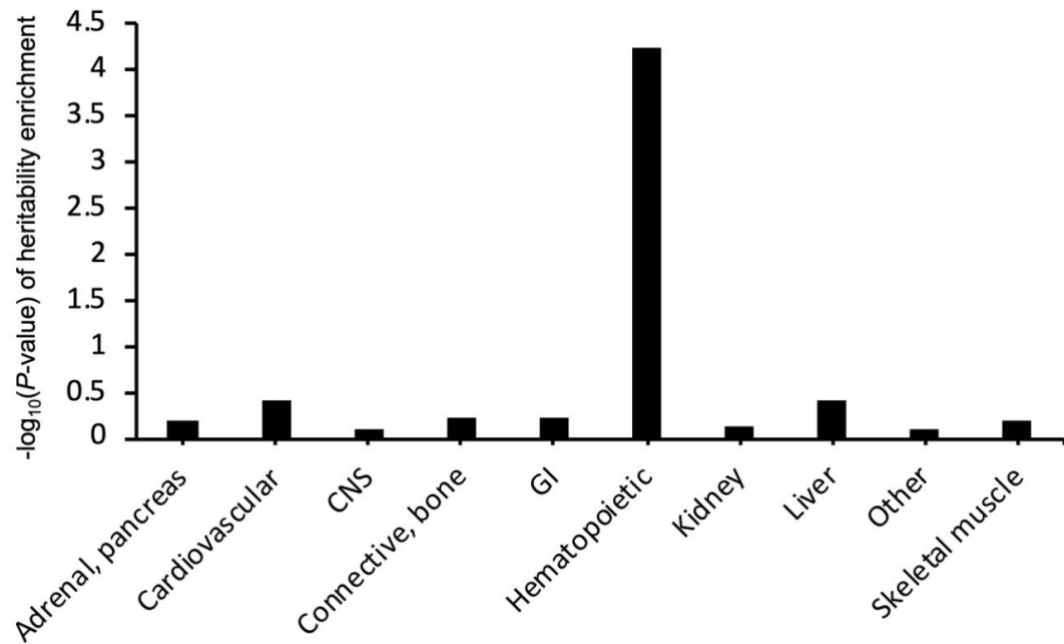


906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924

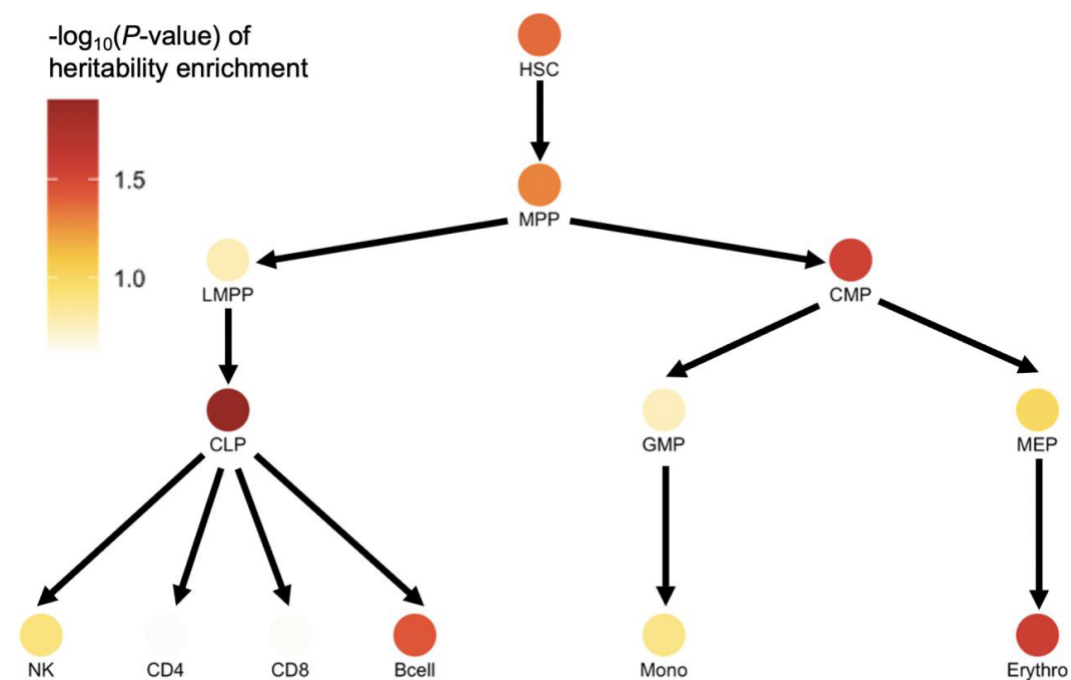
Fig. 2: Associations between CH and diverse traits/diseases.

a-b, Heatmaps showing associations between overall CH, CH with large clones, and CH driven by *DNMT3A*, *TET2*, *ASXL1*, *JAK2*, and *SRSF2+SF3B1* mutations and: **a**, blood cell counts/indices or **b**, biochemical analytes. Colors depict statistical significance of differences compared to individuals without CH, as signed false discovery rate (FDR) values. **c**, Forest plot showing the odds ratios (ORs) for associations between CH and selected traits/diseases prevalent in UKB participants at baseline. **d**, Forest plot showing the ORs for associations between CH subtypes and smoking status, for previous and current smokers. **e-h**, Forest plots showing the hazard ratios (HRs) for associations between CH at baseline and subsequent: **e**, hematological neoplasms, **f**, other malignant neoplasms, **g**, cardiovascular diseases, and **h**, selected causes of death. For **g** and **h**, both overall CH and CH characterized by large clones (“CH large”) are shown. ORs/HR markers with a *P*-value<0.05 are depicted in blue. Error bars represent 95% confidence intervals (CIs). Numerical values for ORs/HRs, 95% CIs, and *P*-values are reported in Supplementary Tables 5–13. Abbreviations: RDW, red blood cell (erythrocyte) distribution width; PDW, platelet distribution width; PCT, plateletcrit; PLT, platelet count; WBC, white blood cell (leukocyte) count; NE, neutrophil count; HLR, high light scatter reticulocyte count; RET, reticulocyte count; MO, monocyte count; MCV, mean corpuscular volume; HT, hematocrit percentage; HGB, hemoglobin concentration; LY, lymphocyte count; RBC, red blood cell (erythrocyte) count; EO, eosinophil count; CYS, cystatin C; PHOS, phosphate; AST, aspartate aminotransferase; HBA1C, glycosylated hemoglobin; APOA, apolipoprotein A; APOB, apolipoprotein B; HDL, HDL cholesterol; LDLD, LDL direct cholesterol; CHOL, total cholesterol; T2DM, type 2 diabetes mellitus; MPN, myeloproliferative neoplasms; MDS, myelodysplastic syndromes; AML, acute myeloid leukemia; CMML, chronic myelomonocytic leukemia.

a



b



925
926
927
928
929
930
931
932
933
934

Fig. 3: Cell type-specific enrichment of the CH polygenic signal.

a, Heritability enrichment of CH across histone marks profiled in 10 cell type groups. **b**, Heritability enrichment of CH across open chromatin regions identified by ATAC-seq in hematopoietic progenitor cells/lineages at different stages of differentiation. Partitioned heritability cell-type group analysis in the LDSC software was used to compute these enrichments and corresponding P -values. The data underlying the figures is available in Supplementary Tables 14 and 15. Abbreviations: CNS, central nervous system; GI, gastrointestinal; CLP, common lymphoid progenitor; CMP, common myeloid progenitor; MPP, multipotent progenitor; HSC, hematopoietic stem cell; GMP, granulocyte/macrophage progenitor; LMPP, lymphoid-primed multipotent progenitor; NK, natural killer cell; Mono, monocyte; Erythro, erythroid progenitor; LDSC, linkage disequilibrium score regression; ATAC-seq, (Assay for Transposase-Accessible Chromatin using sequencing).

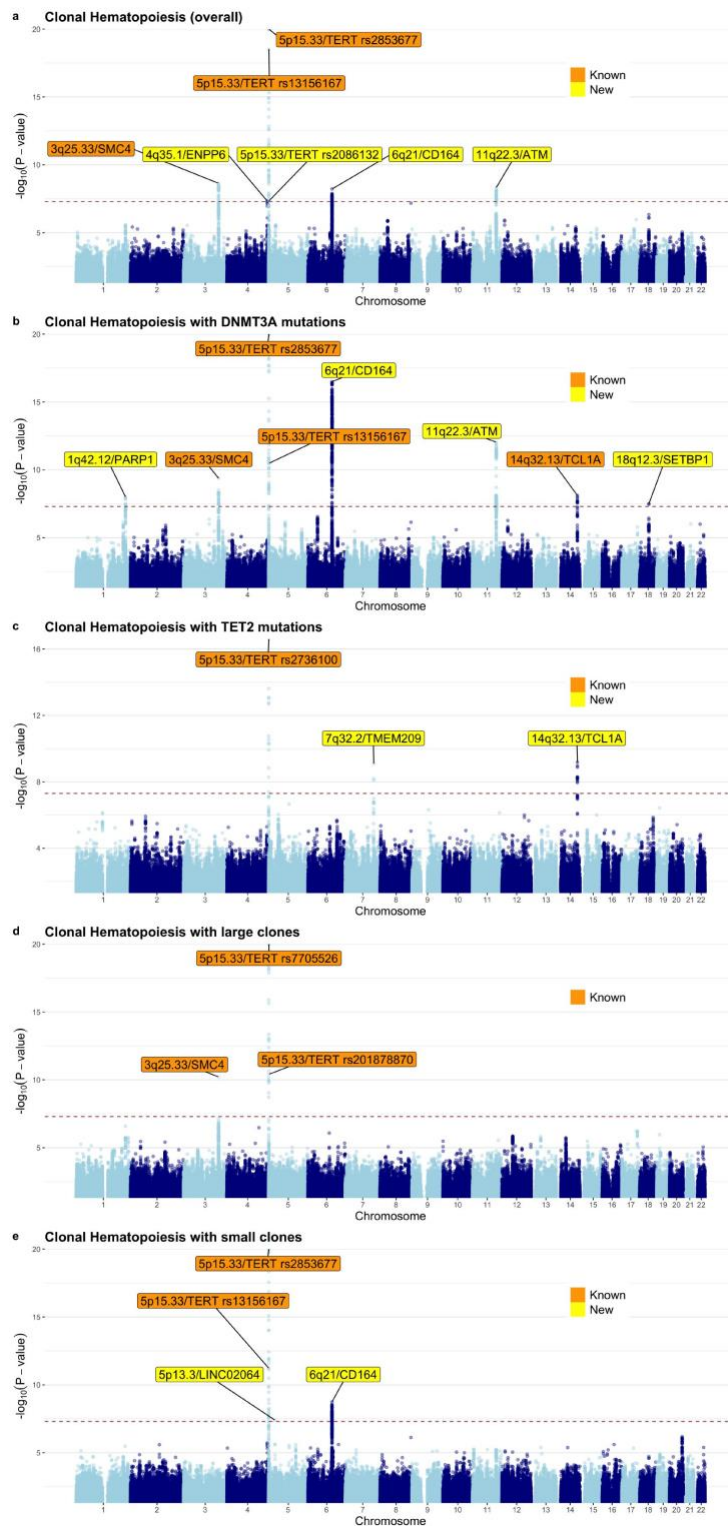


Fig. 4: Manhattan plots displaying genome-wide associations between common germline genetic variants and each of five CH traits.

The Y-axes depict P -values ($-\log_{10}$) for associations derived from the non-infinitesimal mixed model association test implemented in BOLT-LMM. The X-axes depict chromosomal position on build 37 of the human genome (GRCh37). The dotted lines indicate the genome-wide significance threshold of $P=5 \times 10^{-8}$. Known (previously published) and new loci are indicated by cytoband and target gene (based on the prioritization exercise described in the text). Since there were multiple independent loci at 5p15.33 ($LD r^2 < 0.05$), we also label the 5p15.33 signals using the lead variant rs number for each signal. Our prioritization exercise was focused on protein coding genes near each lead variant and since there were no protein coding genes within 1 Mb of the lead variant at 5p13.3, we labeled this association using the nearest non-coding RNA. The CH traits corresponding to each Manhattan plot are: **a**, overall CH. **b**, CH with mutant *DNMT3A*. **c**, CH with mutant *TET2*. **d**, CH with large clones. **e**, CH with small clones.

935
936
937
938
939
940
941
942
943
944
945
946

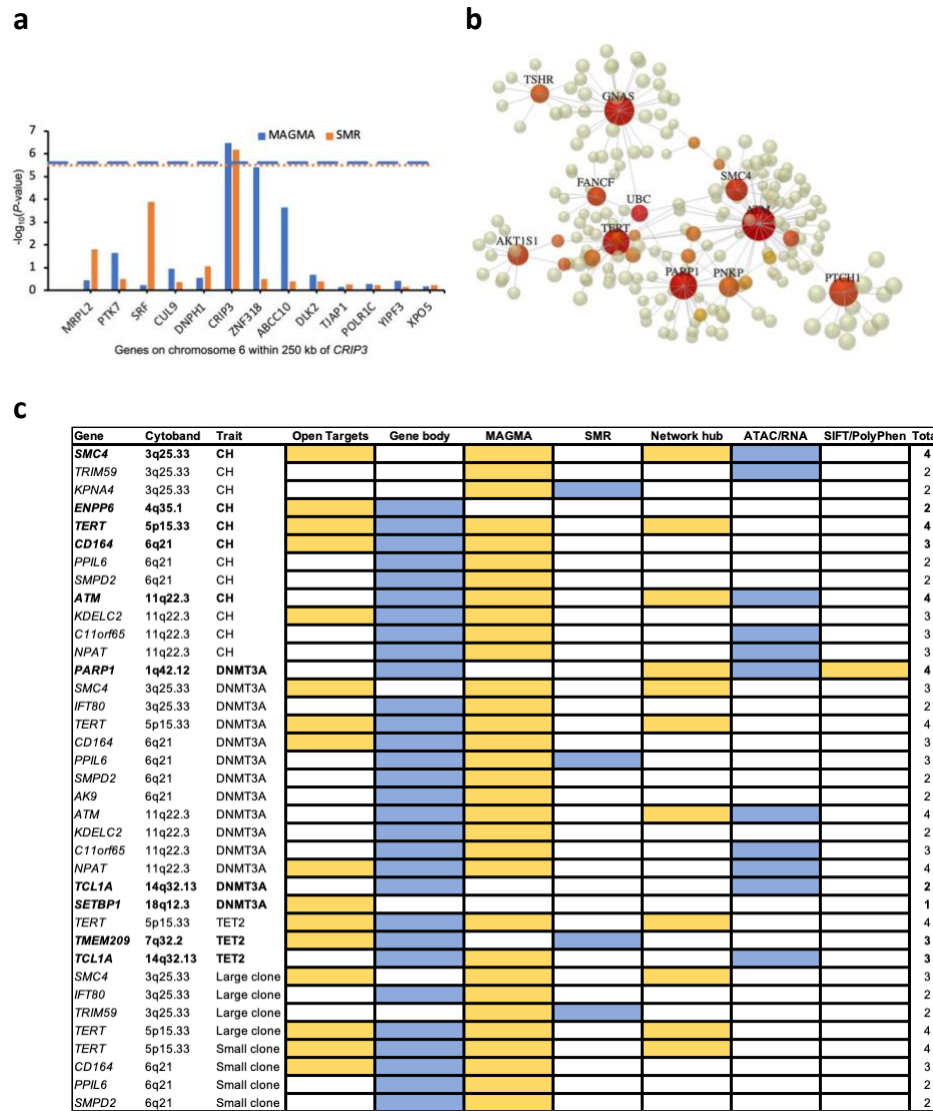
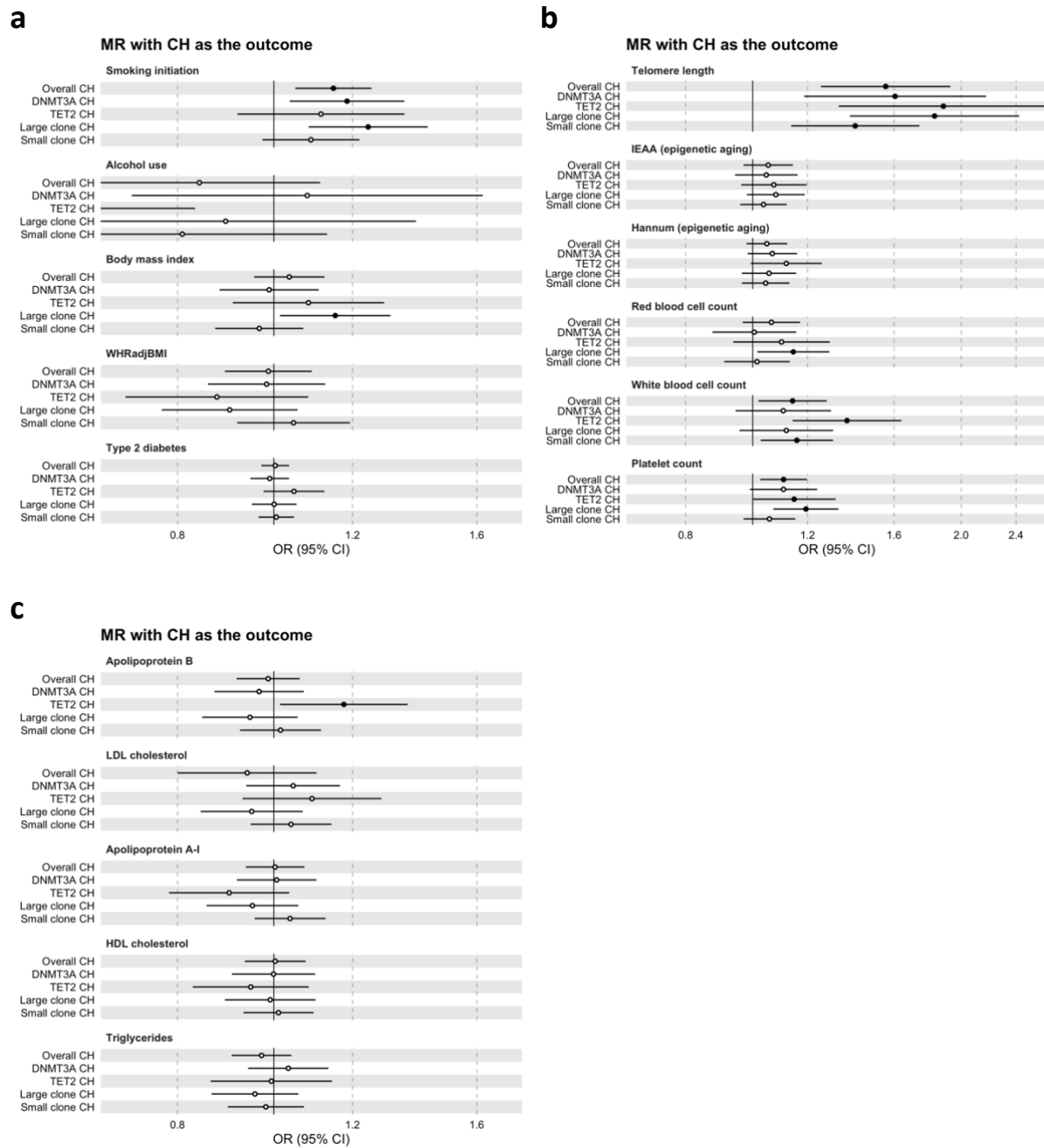


Fig. 5: Gene-level association and protein-protein interaction network analyses, and functional target gene prioritization matrix.

a, Gene-level associations in the 6q21 region within 250 kilobases of *CRIP3*, i.e., between GRCh37 positions 43,017,448 and 43,526,535 on chromosome 6. The X-axis lists all the genes in this region that were tested by both MAGMA and SMR. *CRIP3* was the only gene located more than one megabase away from a GWAS-identified lead variant that was found to be associated with CH at gene-level genome-wide significance by both MAGMA and SMR. The Y-axis depicts the *P*-value ($-\log_{10}$) for association in the MAGMA and SMR analyses. The gene-level genome-wide significance threshold in MAGMA ($P=2.6 \times 10^{-6}$ after accounting for 19,064 genes tested) is indicated by the blue dashed line and in SMR ($P=3.2 \times 10^{-6}$ after accounting for 15,672 genes tested) by the orange dotted line. Both *CRIP3* and *SRF* had SMR HEIDI $P > 0.05$ indicating colocalization of the GWAS and expression quantitative trait locus associations. **b**, Largest sub-network of genes/proteins associated with overall CH risk. All genes ($n=57$) with $P_{\text{MAGMA}} < 0.001$ in the overall CH MAGMA analysis were mapped to proteins and used as “seeds” for network construction that was done by integrating high-confidence protein-protein interactions from the STRING database. The largest sub-network constructed contained 13 of the 57 seed proteins and included 210 nodes and 231 edges. The colored nodes indicate seed proteins that interact with at least two other proteins in this sub-network with the intensity of redness increasing with number of interacting proteins. Seed proteins that interact with six or more other proteins in the sub-network are named above their corresponding node. **c**, Matrix of target genes (protein coding) prioritized by seven approaches (Open Targets, fine-mapped variant-gene body overlap, MAGMA, SMR, network analysis hub, fine-mapped variant-ATAC-seq peak overlap followed by ATAC-RNA-seq correlation, and SIFT/PolyPhen scores) across the loci identified in this study. Only genes prioritized by at least two methods are shown, with the exception of *SETBP1* that is shown despite being prioritized by only one method since it was the only gene prioritized at 18q12.3 and also happens to be an occasional somatic driver of CH. There were no protein coding genes within 1 Mb of the new small clone CH lead variant rs72755524 at 5p13.3 (nearest non-protein coding gene: *LINC02064*, nearest protein coding gene: *CDH6*) and this locus is not shown in the matrix. Abbreviations: MAGMA, multi-marker analysis of genomic annotation; SMR, summary-based Mendelian randomization; HEIDI, heterogeneity in dependent instruments test; SIFT, Sorting Tolerant From Intolerant.

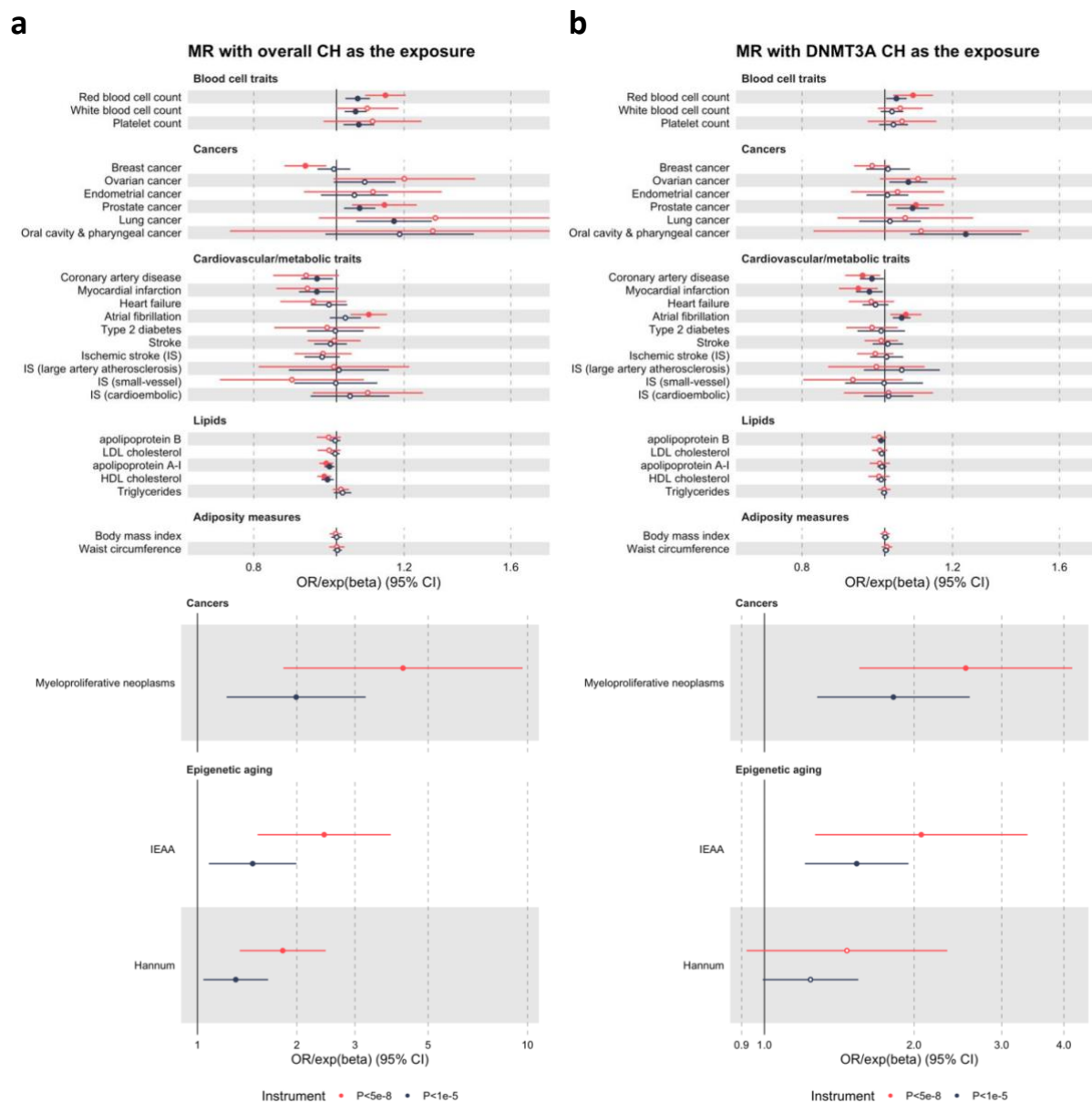
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972



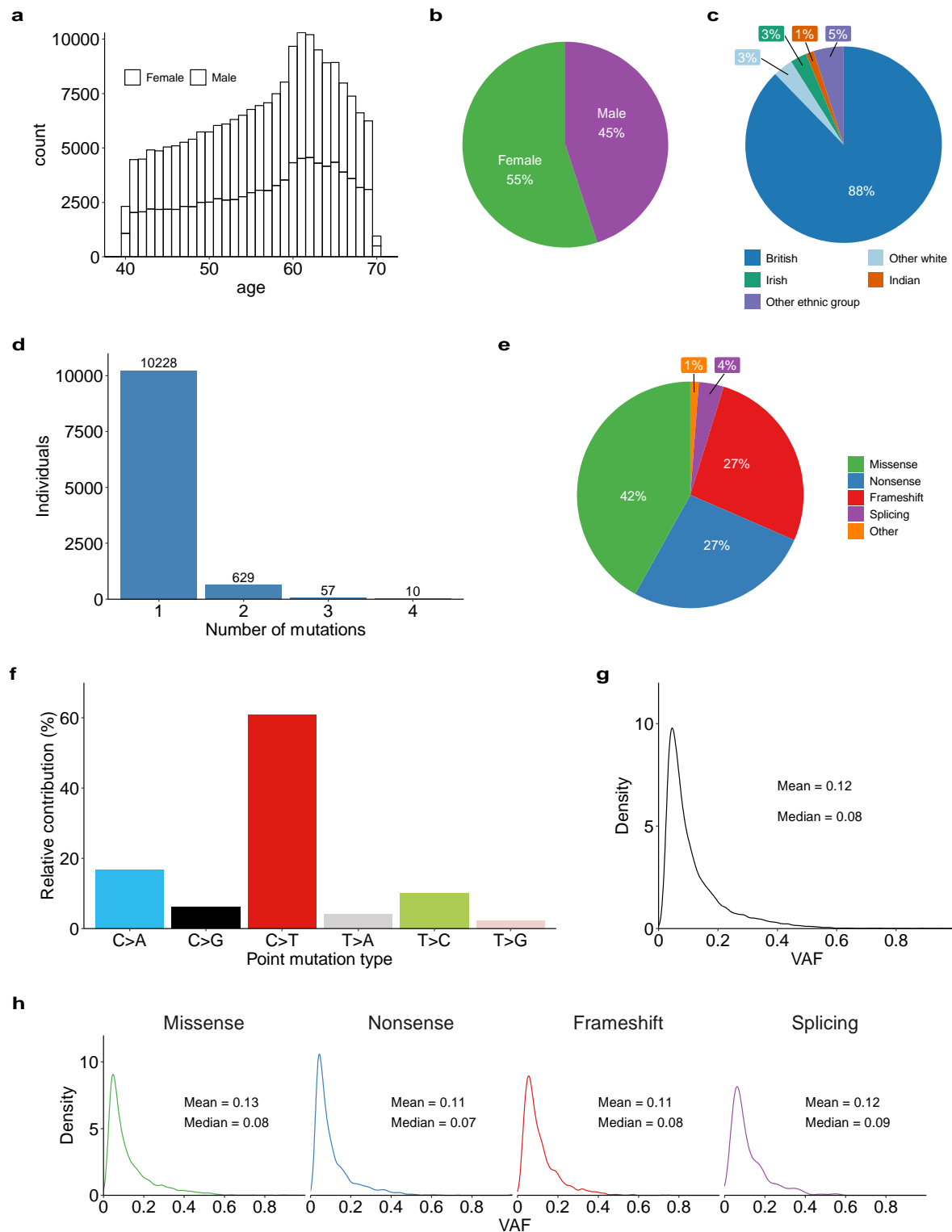
973
974
975
976
977
978
979
980
981
982
983

Fig. 6: Two-sample inverse-variance-weighted Mendelian randomization forest plots with CH traits as outcomes.

Odds ratios (ORs) for CH risk are represented as per (i) standard deviation unit for continuous exposures (alcohol use in drinks per week, body mass index, waist-to-hip ratio adjusted for body mass index (WHRadjBMI) (a); leukocyte telomere length, two epigenetic aging traits, red cell, white cell and platelet counts (b), and five circulating lipid traits (c) and (ii) log-odds unit for binary exposures (smoking initiation (ever having smoked regularly) and genetic liability to type 2 diabetes (a)). Details of units are provided in Supplementary Table 30. OR markers with corresponding P -value < 0.05 are represented by filled circles. Error bars represent 95% confidence intervals (CIs). Full results, including sensitivity analyses, are presented in Supplementary Tables 32, 33, and 34. Abbreviations: MR, Mendelian randomization; CH, clonal hematopoiesis; WHRadjBMI, waist-to-hip ratio adjusted for body mass index; LDL, low-density lipoprotein cholesterol; HDL, high-density lipoprotein cholesterol; IEAA, intrinsic epigenetic age acceleration.



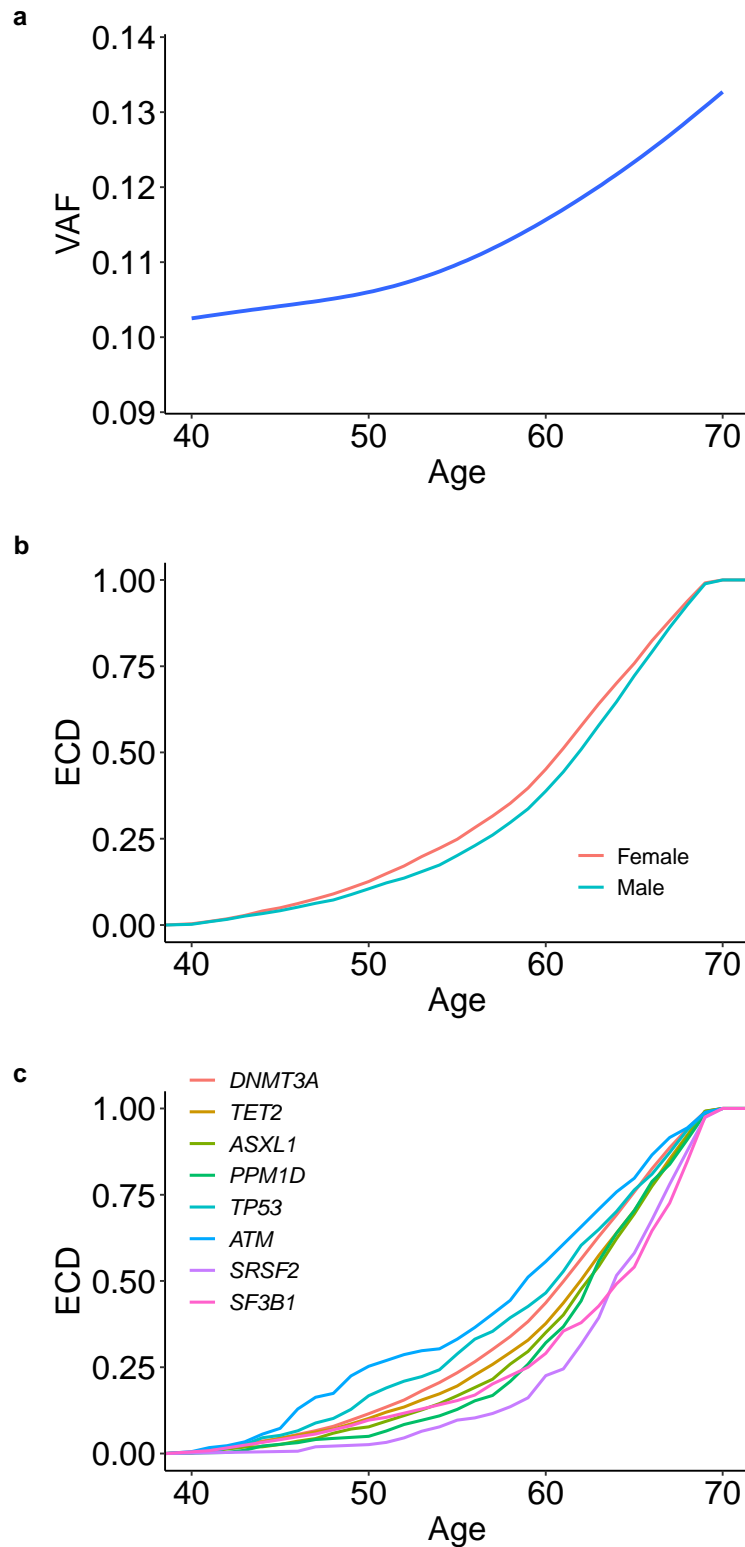
984
985 **Fig. 7: Two-sample inverse-variance-weighted Mendelian randomization forest plots with CH traits as exposures.**
986 Forest plots with odds ratio (OR) markers (for cancers and cardiovascular/metabolic traits) or exponentiated beta coefficient
987 (exp(beta)) markers (for blood cell traits, lipids, adiposity measures, and epigenetic aging indices). ORs/exp(betas) are
988 represented as per log-odds unit increase in genetic liability to **a**, overall CH, or **b**, DNMT3A CH. OR/exp(beta) markers with
989 corresponding P -value<0.05 are represented by filled circles. Error bars represent 95% confidence intervals (CIs). Red
990 markers and error bars represent results using genetic instruments comprised exclusively of genome-wide significant
991 ($P < 5 \times 10^{-8}$) variants. Black markers and error bars represent results when using genome-wide significant and sub-genome-
992 wide significant ($P < 10^{-5}$) variants in the genetic instrument. Large effect size estimates (ORs/exp(betas)) are shown in the
993 lower panels. Full results, including sensitivity analyses, are presented in Supplementary Tables 35 and 36. Abbreviations:
994 MR, Mendelian randomization; IS, ischemic stroke; LDL, low-density lipoprotein cholesterol; HDL, high-density lipoprotein
995 cholesterol; IEAA, intrinsic epigenetic age acceleration.



996
997
998
999
1000
1001
1002
1003
1004
1005

Extended Data Fig. 1: Characterization of CH in the UK Biobank.

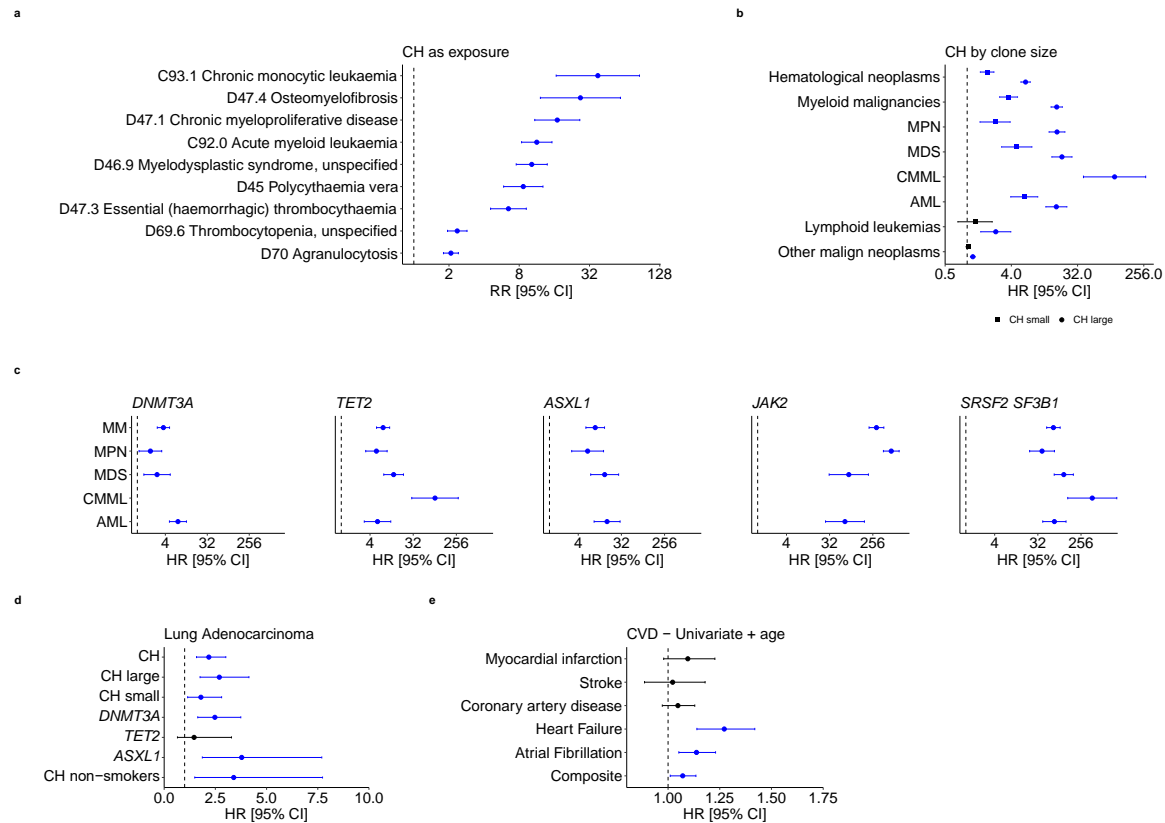
a, Histogram stratified by sex showing the age-distribution of individuals in the UKB cohort (n=200,453). **b**, Overall percentage of females and males in the UKB cohort. **c**, Percentage of the most common self-reported ethnic groups in the UKB cohort. Ethnic groups with a frequency lower than 1% were grouped under the “Other ethnic group” category. **d**, Number of individuals with 1, 2, 3, and 4 somatic mutations. More than 90% of individuals with CH had only one driver mutation identified. **e**, Percentages of different CH mutation types identified. **f**, Relative prevalence of each of the six base substitution types amongst the identified CH mutations. **g**, Density plot showing the variant allele fraction (VAF) distribution of all CH somatic mutations. **h**, Density plot showing similar VAF distribution for different mutation types. Mean and median are indicated for **g** and **h**.



1006
1007
1008
1009
1010
1011
1012
1013
1014
1015

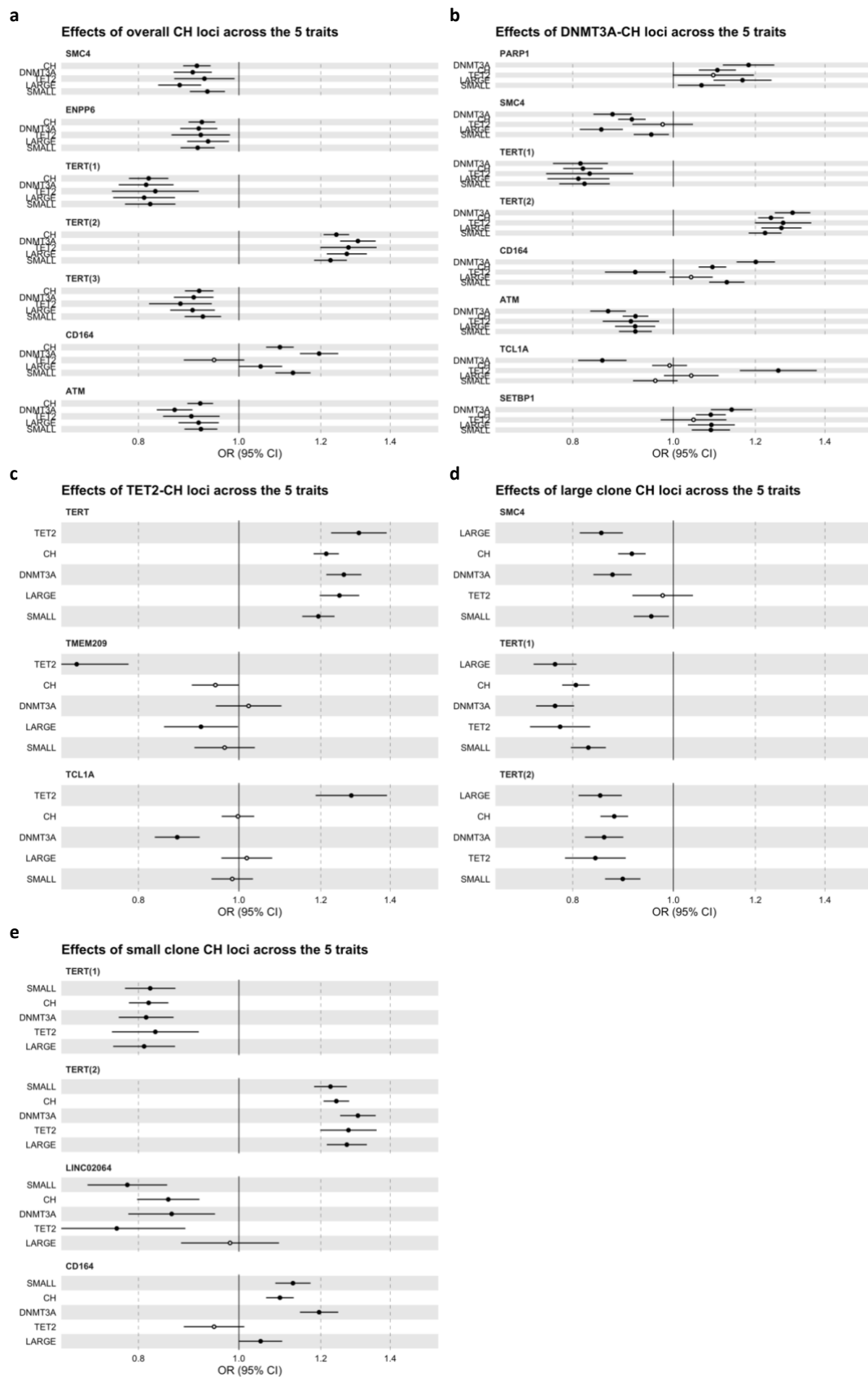
Extended Data Fig. 2: Age-distribution of CH by clone size, sex, and mutant gene.

a, Clone size, estimated by the variant allele fraction (VAF), increases with age. The blue line represents the smoothed model fitted to a generalized additive model and the shadow represents the 95% confidence interval. **b**, Empirical cumulative distribution (ECD) of the age of individuals with CH stratified by sex. CH was observed one year earlier in females than in males (median 61 versus 62 years; $P=1.6 \times 10^{-4}$, two-sided pairwise Wilcoxon rank sum test). **c**, ECD of the age of individuals with CH stratified by the eight most common driver genes. Compared to *DNMT3A*, mutations in *ATM* were observed 3 years earlier ($P=7.2 \times 10^{-4}$), while mutations in *ASXL1*, *PPM1D*, *SRSF2*, and *SF3B1* were observed 1 ($P=2.7 \times 10^{-8}$), 1 ($P=8.5 \times 10^{-6}$), 2 ($P=5.7 \times 10^{-10}$), and 3 ($P=6.5 \times 10^{-6}$) years later, respectively. Differences were calculated using a pairwise Wilcoxon rank sum test.



Extended Data Fig. 3: Associations between CH and diseases.

a, Association analysis of CH with International Classification of Diseases version-10 (ICD-10) disease codes, showing the risk ratios (RRs) for ICD-10 codes with CH as exposure. Only ICD-10 codes with false discovery rate (FDR) $<10^{-10}$ are represented. Error bars represent 95% confidence intervals (CIs). **b-c**, Forest plots showing the hazard ratios (HRs) from Cox proportional-hazards models for association with subsequent hematological and other malignant neoplasms for CH with small and large clones (**b**) and CH driven by mutations in specific genes (**c**). **d**, Forest plot showing the HRs for subsequent/incident lung adenocarcinoma for overall CH, CH with large and small clones, and CH with *DNMT3A*, *TET2*, and *ASXL1* mutations, as well as for overall CH restricting to only self-reported “never-smokers”. **e**, HRs for subsequent/incident cardiovascular disease (CVD) conditions after CH at baseline using a bivariable model containing age as the only covariate. For **b-e**, HR markers with P -value <0.05 are depicted in blue. Error bars represent 95% CIs. Numerical values for RRs/HRs, 95% CIs, and P -values are reported in Supplementary Tables 9–12.



Extended Data Fig. 4: Heterogeneity of lead GWAS variants across five CH traits.

Forest plots with odds ratios (ORs) and 95% confidence intervals (CIs) based on data from Supplementary Tables a, 16, b, 18, c, 19, d, 20, and e, 21. Results for lead variants identified at genome-wide significance ($P < 5 \times 10^{-8}$) for each CH trait (a, overall CH, b *DNMT3A*-CH, c *TET2*-CH, d large clone CH, and e, small clone CH) are plotted alongside results for the same lead variants in the four other genome-wide association analyses conducted.

1028
1029
1030
1031
1032
1033