

# 1 The Genetic Determinants and Genomic Consequences 2 of Non-Leukemogenic Somatic Point Mutations 3

4 Authors: Joshua S. Weinstock<sup>1,#</sup>, Sharjeel A. Chaudhry<sup>2,3</sup>, Maria Ioannou<sup>4</sup>, Maria Viskadourou<sup>2</sup>, Paula  
5 Reventun<sup>2</sup>, Yasminka A. Jakubek<sup>5</sup>, L. Alexander Liggett<sup>6</sup>, Cecelia Laurie<sup>7</sup>, Jai G. Broome<sup>8</sup>, Alyna Khan<sup>9</sup>,  
6 Kent D. Taylor<sup>10</sup>, Xiuqing Guo<sup>10</sup>, Patricia A. Peyser<sup>11</sup>, Eric Boerwinkle<sup>12</sup>, Nathalie Chami<sup>13,14</sup>, Eimear E.  
7 Kenny<sup>15</sup>, Ruth J. Loos<sup>13,14</sup>, Bruce M. Psaty<sup>16,17,18</sup>, Tracy P. Russell<sup>19</sup>, Jennifer A. Brody<sup>16</sup>, Jeong H. Yun<sup>20</sup>,  
8 Michael H. Cho<sup>21</sup>, Ramachandran S. Vasan<sup>22</sup>, Sharon L. Kardia<sup>23</sup>, Jennifer A. Smith<sup>23,24</sup>, Laura M. Raffield<sup>25</sup>,  
9 Aurelian Bidulescu<sup>26</sup>, Emily O'Brien<sup>27</sup>, Mariza de Andrade<sup>28</sup>, Jerome I. Rotter<sup>10</sup>, Stephen S. Rich<sup>29</sup>, Russell  
10 P. Tracy<sup>19</sup>, Yii Der Ida Chen<sup>10</sup>, C. Charles. Gu<sup>30</sup>, Chao A. Hsiung<sup>31</sup>, Charles Kooperberg<sup>32</sup>, Bernhard  
11 Haring<sup>33,34</sup>, Rami Nassir<sup>35</sup>, Rasika Mathias<sup>36</sup>, Alex Reiner<sup>32</sup>, Vijay Sankaran<sup>6</sup>, Charles J. Lowenstein<sup>37</sup>,  
12 Thomas W. Blackwell<sup>38</sup>, Goncalo R. Abecasis<sup>38,39</sup>, Albert V. Smith<sup>38</sup>, Hyun M. Kang<sup>38</sup>, Pradeep  
13 Natarajan<sup>40,41,42</sup>, Siddhartha Jaiswal<sup>43</sup>, Alexander Bick<sup>44</sup>, Wendy S. Post<sup>37</sup>, Paul Scheet<sup>45</sup>, Paul Auer<sup>46</sup>,  
14 Theodoros Karantanos<sup>4</sup>, Alexis Battle<sup>47,48,49,#</sup>, Marios Arvanitis<sup>2,47,#</sup>

15  
16 1 - Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, USA; 2 - Division  
17 of Cardiology, Department of Medicine, Johns Hopkins University, Baltimore, MD; 3 - Department of  
18 Surgery, Division of Vascular and Endovascular Surgery, Beth Israel Deaconess Medical Center, Harvard  
19 Medical School, Boston, MA, USA; 4 - Division of Hematological Malignancies, Department of Oncology,  
20 Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins University School of Medicine; 5 -  
21 Department of Internal Medicine, University of Kentucky; 6 - Division of Hematology/Oncology, Boston  
22 Childrens Hospital and Department of Pediatric Oncology, Dana-Farber Cancer Institute, Harvard  
23 Medical School, Boston, MA 02115, USA; 7 - Department of Biostatistics, University of Washington,  
24 Seattle, WA 98195, USA.; 8 - Division of Medical Genetics, Department of Medicine, University of  
25 Washington, Seattle, WA 98195, USA; 9 - Department of Biostatistics, University of Washington, Seattle,  
26 WA 98195, USA; 10 - The Institute for Translational Genomics and Population Sciences, Department of  
27 Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance,  
28 CA USA; 11 - Department of Epidemiology, School of Public Health, Boston University, Boxton, MA USA;  
29 12 - Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA; 13 - The Charles  
30 Bronfman Institute of Personalized Medicine; 14 - The Mindich Child Health and Development  
31 Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA; 15 - Institute for Genomic  
32 Health; 16 - Cardiovascular Health Research Unit, Department of Medicine, University of Washington,  
33 Seattle, WA, USA; 17 - Department of Epidemiology, University of Washington, Seattle, WA, USA; 18 -  
34 Department of Health Systems and Population Health, University of Washington, Seattle, WA, USA; 19 -  
35 Department of Pathology & Laboratory Medicine and Biochemistry, Larner College of Medicine at the  
36 University of Vermont, Colchester, VT, USA; 20 - Channing Division of Network Medicine, Brigham and  
37 Women's Hospital, Boston, MA USA; 21 - Channing Division of Network Medicine and Division of  
38 Pulmonary and Critical Care Medicine, Brigham and Women's Hospital, Boston, MA USA; 22 - National  
39 Heart Lung and Blood Institute's, Boston University's Framingham Heart Study, Framingham, MA, USA;  
40 23 - Department of Epidemiology, University of Michigan, Ann Arbor, MI; 24 - Survey Research Center,  
41 Institute for Social Research, University of Michigan, Ann Arbor, MI; 25 - Department of Genetics,  
42 University of North Carolina, Chapel Hill, NC, 27514; 26 - Department of Epidemiology and Biostatistics,  
43 Indiana University School of Public Health Bloomington, Bloomington, IN, USA; 27 - Duke Clinical  
44 Research Institute, Durham, NC, USA; 28 - Mayo Clinic, Department of Health Sciences Research,  
45 Rochester, MN, USA; 29 - Department of Public Health Sciences, Center for Public Health Genomics,

46 University of Virginia, Charlottesville, VA USA; 30 - Center for Biostatistics and Data Sciences,  
47 Washington University, St. Louis, MO USA; 31 - Department of Medicine, Taipei Veterans General  
48 Hospital, Taipei Taiwan; 201 Shi-Pai Rd. Sec. 2, Taipei Taiwan; 32 - Division of Public Health Sciences,  
49 Fred Hutchinson Cancer Research Center, Seattle, WA, USA; 33 - Department of Medicine III, Saarland  
50 University Hospital, Homburg, Saarland, Germany; Department of Medicine I, University of Würzburg,  
51 Würzburg, Bavaria, Germany; 34 - Department of Epidemiology and Population Health, Albert Einstein  
52 College of Medicine, Bronx, New York, USA. Electronic address; 35 - University of California Davis, Davis,  
53 CA, USA; 36 - Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD,  
54 USA; 37 - Department of Medicine, Cardiology Division, Johns Hopkins University; 38 - Center for  
55 Statistical Genetics, Department of Biostatistics, University of Michigan School of Public Health, Ann  
56 Arbor, MI, USA; 39 - Regeneron Pharmaceuticals, Tarrytown, NY, USA; 40 - Center for Genomic Medicine  
57 and Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA; 41 - Program in  
58 Medical and Population Genetics, Broad Institute of Harvard & MIT, Cambridge, MA; 42 - Department of  
59 Medicine, Harvard Medical School, Boston, MA; 43 - Department of Pathology, Stanford University,  
60 Stanford, CA, USA; 44 - Division of Genetic Medicine, Department of Medicine, Vanderbilt University,  
61 Nashville, TN, USA; 45 - Department of Epidemiology, University of Texas M.D. Anderson Cancer Center,  
62 Houston, TX, USA.; 46 - Department of Biostatistics, Medical College of Wisconsin Division of  
63 Biostatistics, Institute for Health and Equity, and Cancer Center, Medical College of Wisconsin,  
64 Milwaukee, WI, USA; 47 - Department of Biomedical Engineering, Johns Hopkins University, Baltimore,  
65 MD, USA; 48 - Malone Center for Engineering in Healthcare, Johns Hopkins University, Baltimore, MD;  
66 49 - Department of Computer Science, Johns Hopkins University, Baltimore, MD

67  
68 # Address correspondence to:

69 Marios Arvanitis (marvani1@jhmi.edu), Josh Weinstock (josh.weinstock@emory.edu), Alexis Battle  
70 (ajbattle@jhu.edu)

71

72 Abstract

73 Clonal hematopoiesis (CH) is defined by the expansion of a lineage of genetically identical cells in blood.  
74 Genetic lesions that confer a fitness advantage, such as point mutations or mosaic chromosomal  
75 alterations (mCAs) in genes associated with hematologic malignancy, are frequent mediators of CH.  
76 However, recent analyses of both single cell-derived colonies of hematopoietic cells and population  
77 sequencing cohorts have revealed CH frequently occurs in the absence of known driver genetic lesions.  
78 To characterize CH without known driver genetic lesions, we used 51,399 deeply sequenced whole  
79 genomes from the NHLBI TOPMed sequencing initiative to perform simultaneous germline and somatic  
80 mutation analyses among individuals without leukemogenic point mutations (LPM), which we term CH-  
81 LPMneg. We quantified CH by estimating the total mutation burden. Because estimating somatic  
82 mutation burden without a paired-tissue sample is challenging, we developed a novel statistical method,  
83 the Genomic and Epigenomic informed Mutation (GEM) rate, that uses external genomic and  
84 epigenomic data sources to distinguish artifactual signals from true somatic mutations. We performed a  
85 genome-wide association study of GEM to discover the germline determinants of CH-LPMneg. After  
86 fine-mapping and variant-to-gene analyses, we identified seven genes associated with CH-LPMneg  
87 (*TCL1A*, *TERT*, *SMC4*, *NRIP1*, *PRDM16*, *MSRA*, *SCARB1*), and one locus associated with a sex-associated  
88 mutation pathway (*SRGAP2C*). We performed a secondary analysis excluding individuals with mCAs,  
89 finding that the genetic architecture was largely unaffected by their inclusion. Functional analyses of  
90 *SMC4* and *NRIP1* implicated altered HSC self-renewal and proliferation as the primary mediator of

91 mutation burden in blood. We then performed comprehensive multi-tissue transcriptomic analyses,  
92 finding that the expression levels of 404 genes are associated with GEM. Finally, we performed  
93 phenotypic association meta-analyses across four cohorts, finding that GEM is associated with increased  
94 white blood cell count and increased risk for incident peripheral artery disease, but is not significantly  
95 associated with incident stroke or coronary disease events. Overall, we develop GEM for quantifying  
96 mutation burden from WGS without a paired-tissue sample and use GEM to discover the genetic,  
97 genomic, and phenotypic correlates of CH-LPMneg.

98

## 99 Introduction

100 As we age, our cells accumulate mutations. The vast majority of these mutations are  
101 inconsequential because they do not alter cell fitness. However, a small proportion of these mutations,  
102 termed drivers, can cause expansions of cell lineages they reside in. Recently, the age-related acquisition  
103 of leukemogenic point mutations (LPM) in whole blood, termed clonal hematopoiesis of indeterminate  
104 potential (CHIP), has been described as a prevalent aging-related phenomenon<sup>1-4</sup>. CHIP has previously  
105 been associated with increased risk for hematologic malignancy, cardiovascular disease, and increased  
106 mortality<sup>4-6</sup>. However, CHIP is a highly specific clonal phenomena defined as the presence of a driver  
107 mutation in 74 genes that have previously been associated with hematologic malignancy<sup>7</sup>, which is a  
108 small proportion of the entire spectrum of somatic variation. Non-CHIP somatic variation in blood, which  
109 we term CH-LPMneg, has previously been shown to be a prevalent phenomenon, including mosaic  
110 chromosomal alterations (mCAs)<sup>8-10</sup>, X-chromosome inactivation skewing<sup>11</sup>, and even clonal expansions  
111 without known drivers<sup>11,12</sup>. However, the germline determinants and clinical consequences of CH-  
112 LPMneg remain uncharacterized.

113 We previously used the count of high variant-allele fraction (VAF) passenger mutations in 5,071  
114 CHIP carriers in TOPMed to identify the genetic determinants of clonal expansion, an approach termed  
115 PACER<sup>13</sup>, which uses age at blood-draw and passenger burden to infer the date at which a driver  
116 mutation was acquired. However, PACER is only defined for donors with a single driver mutation. Here,  
117 we seek to extend our inference of the sample level mutation burden for donors that may have no  
118 known driver point mutations. Non-CHIP clonal phenomena, which we refer to as CH-LPMneg, including  
119 mCAs, LOY, X-chromosome inactivation, and clonal expansions without known drivers have been  
120 previously associated with infection<sup>14</sup>, hematologic malignancy<sup>11</sup>, and heart failure<sup>15</sup>, highlighting the  
121 value of quantifying CH-LPMneg.

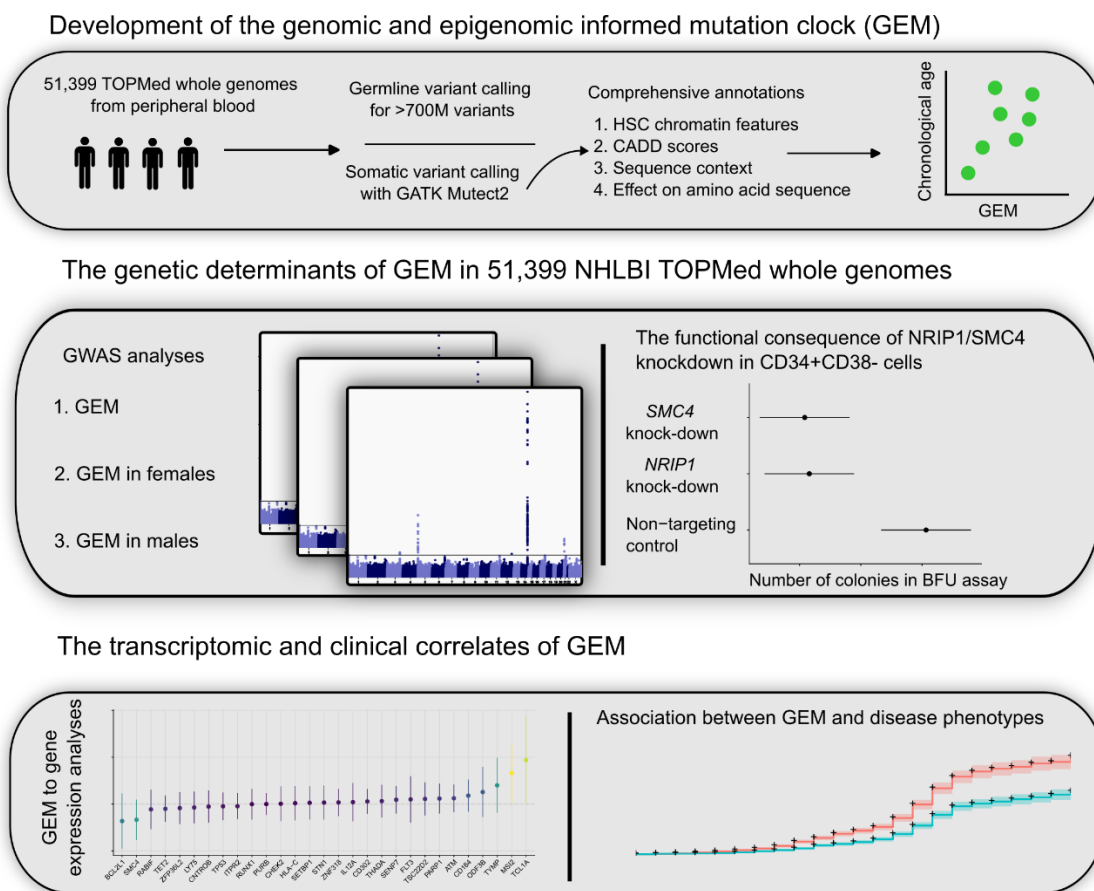
122 The accurate detection of somatic mutations in CHIP non-carriers from a single whole-blood  
123 draw is likely to be more challenging than the identification in CHIP carriers because the passenger  
124 count no longer tracks the history of a single expanded clone. We reasoned that improved estimation of  
125 the somatic mutation rate would facilitate more accurate passenger burden inference in this more  
126 challenging setting. Previous reports have identified chromatin state as among the primary  
127 determinants of mutation rate. Indeed, Shuster-Bockler and Lehner<sup>16</sup> reported that variation in  
128 chromatin organization explains 55% of the variation in mutation rate. Thus, external epigenomic  
129 annotations are informative for estimating the likelihood that a candidate somatic variant call is a true  
130 mutation by altering the prior probability that a given variant call is accurate.

131 We and others have previously reported that CHIP and other clonal phenomena have germline  
132 genetic determinants<sup>9,11,13,17-20</sup>. CHIP has been previously associated with two primary pathways  
133 influencing HSC self-renewal and DNA damage pathways. Similarly, GWAS of mCAs have similarly  
134 reported associations with HSC self-renewal and DNA repair related loci<sup>21-23</sup>. A recent analysis that  
135 defined CH based on the dichotomization of low-VAF mutation burden, termed barcode-CH<sup>24</sup>, observed  
136 several hits linked to both mCAs and CHIP. These analyses have demonstrated that germline variation is  
137 associated with acquired genetic variation and have demonstrated the utility of such analyses for  
138 discovering critical regulators of clonal expansion rate, including *TCL1A*.

139 Here, using 51,399 donors from NHLBI TOPMed consortium<sup>25</sup> without CHIP, we developed a  
140 mutation burden estimator, the Genomic and Epigenomic Mutation (GEM) rate (Figure 1). We then used  
141 this estimator as a phenotype to discover the genetic determinants of CH-LPMneg. In contrast to  
142 barcode-CH, this is a continuous phenotype that excludes individuals with CHIP mutations. This analysis  
143 revealed multiple novel loci, including the previously underappreciated role of *NRIP1*, a highly conserved

144 transcriptional co-activator, in modulating mutation burden. We performed a sensitivity analysis,  
 145 excluding all individuals with mCAs, finding that the bulk of our genetic discovery was unchanged,  
 146 suggesting that CH-LPMneg signals are not merely mediated through mCAs. Fine-mapping of the *TCL1A*  
 147 locus revealed a more complex cis-regulatory architecture than observed in CHIP. Functional  
 148 characterization of *SMC4* and *NRIP1* with colony forming unit (CFU) assays revealed convergent effects  
 149 on HSC self-renewal as the primary mechanism of these genes. Sex-stratified analyses of GEM revealed  
 150 that the *TRIM59-KPNA4-SMC4* locus, which has previously been associated with CHIP and MPNs<sup>17,26,27</sup> is  
 151 a female-specific signal, and a novel male-specific signal near *MSRA*. Principal component analysis of  
 152 mutation burden revealed a sex specific mutation pathway. GWAS of this sex-specific mutation pathway  
 153 identified a novel locus near *SRGA2PC*. Through transcriptomic analyses of blood and non-blood tissues,  
 154 we identified the genomic consequences of elevated mutation burden in whole blood, which include the  
 155 systematic down-regulation of the interferon-alpha pathway across hematopoietic lineages. Finally, we  
 156 show that GEM is useful for predicting risk of incident peripheral artery disease, and associates with  
 157 altered blood cell indices. Overall, we demonstrate a novel computational approach for quantifying  
 158 mutation burden, which enabled the discovery of novel genetic determinants of CH-LPMneg.

159 *Figure 1: Study design schematic, describing the development of GEM, the use of GEM to discover the genetic determinants of*  
 160 *mutation burden in blood, and the use of GEM to identify the transcriptomic and clinical correlates of mutation burden in blood.*



161

162 **Results**

163 Using 51,399 WGS samples from NHLBI TOPMed (Supplementary Tables 1-2), we first called  
 164 candidate somatic variants using Mutect2 as previously described<sup>17</sup>. We then performed stringent  
 165 filtering, including filtering known germline variants and likely sequencing artifacts (Methods). As

166 distinguishing somatic variants from germline variants in single-tissue variant calling is challenging, we  
167 then took careful measures to determine the optimal alt-allele threshold. We observed that excluding  
168 variants with higher alt-alleles substantially improved the association of the burden of such mutations  
169 with chronological age, suggesting that a stricter alt-allele threshold than we previously applied in the  
170 PACER pipeline is useful for excluding germline variation.

## 171 Genomic and Epigenomic Annotations Inform Mutation Rate

172

173 Next, as chromatin state is among the primary mediators of mutation rate<sup>16</sup>, we sought to  
174 determine the association between mutation burden and several genomic and epigenomic annotations.  
175 We used chromHMM<sup>28</sup> annotations in CD34+ cells from the Roadmap Epigenomic<sup>29</sup> resource as a  
176 measure of chromatin state in HSCs, which previous analyses have reported as the causal cell type in  
177 clonal phenomena<sup>27</sup>. We calculated the mutation burden stratified by chromatin annotation and  
178 examined the association with age, reasoning that the strength of association between mutation burden  
179 and chronological age would reflect the proportion of artifactual mutations. We observed that  
180 mutations in quiescent chromatin (Figure 2A) are much more strongly associated with age than  
181 mutations in transcriptionally active chromatin, recapitulating the role of chromatin in modifying  
182 mutation rate.

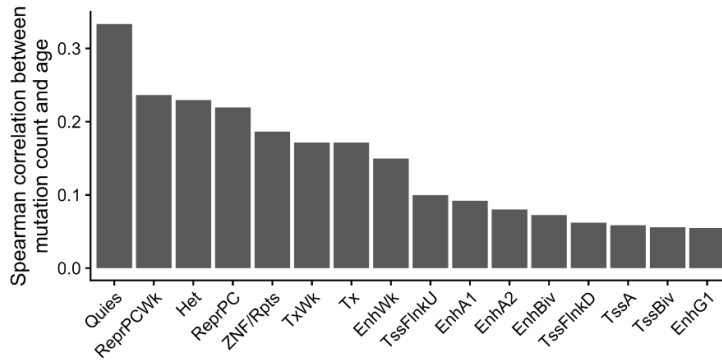
183 As mutations that are functional and are not mutated at the stem cell level undergo extensive  
184 negative selection<sup>30,31</sup>, we then asked whether mutation burden stratified by functional consequence on  
185 protein coding sequence modified the association with chronological age. We observed that mutation  
186 burden of missense and UTR variants was much more weakly associated with chronological age than  
187 mutation burden from intronic and intergenic mutations (Figure 2B), lending credence to our hypothesis  
188 that functional consequence is informative for refining mutation burden estimates. We performed a  
189 similar analysis with stratified CADD scores, finding that the mutations in the highest quintile of CADD  
190 scores were the most weakly associated with age (Figure 2C), again suggesting that deleterious  
191 mutations are depleted of association with age and likely enriched for false positives. We performed an  
192 analyses stratified by both CADD quantile and chromHMM annotation, finding that the two were not  
193 redundant (Extended Data Figure 1). Collectively, these analyses suggest that variant deleteriousness  
194 and chromatin annotations are useful for the construction of a mutation derived molecular clock.

195 We then developed a weakly-supervised probabilistic graphical modeling approach that  
196 incorporates genomic and epigenomic annotations to distinguish somatic mutations from artifacts, a  
197 method we term GEM (genomic and epigenomic mutation rate). Weakly-supervised probabilistic  
198 graphical modeling approaches have been previously used to identify functional rare-variants<sup>32</sup>,  
199 demonstrating the utility of such approaches towards annotation of genetic variation. We first  
200 comprehensively annotated all candidate somatic variants based on their chromHMM<sup>28</sup> annotations,  
201 their functional consequence, CADD<sup>33</sup>, their population allele frequency in TOPMed, surrounding  
202 sequence context, among others (Methods). GEM uses chronological age as an external annotation to  
203 identify which candidate somatic mutations were functional based on their annotations (Figure 2D).  
204 GEM enables the classification of candidate somatic mutations as either truly somatic or artifacts, thus  
205 increasing power in downstream analyses by facilitating the depletion of likely artifactual variants.

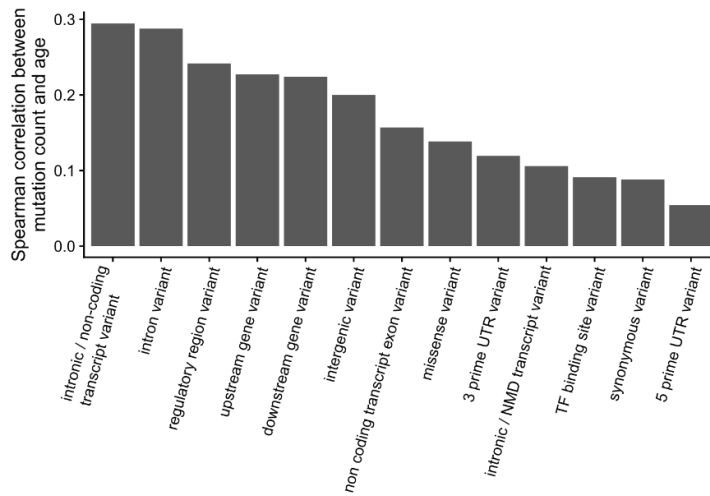
206 *Figure 2: Development of GEM | A, The Spearman correlation between mutation burden and chronological age stratified by*  
207 *chromHMM annotations in CD34+ cells. B, The Spearman correlation between mutation burden and chronological age stratified*  
208 *by functional consequence as annotated by the variant effect predictor (VEP). C, The Spearman correlation between mutation*  
209 *burden and chronological age, stratified by quintiles of CADD scores. D, Plate annotation for the GEM statistical model.  $\theta_0$  and*

210 are intercepts;  $\theta_1$  reflects the association between  $\log_2$  transformed value of  $\sum z_{ij}$  and chronological age  $Y_i$ ;  $z_{ij}$  denotes the  
 211 probability that the  $j$ th mutation in the  $i$ th individual is a true somatic mutation.  $X$  is a matrix of annotations.

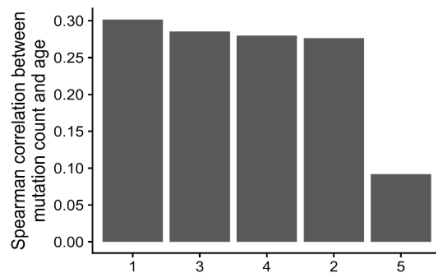
A Mutations in quiescent CD34+ chromatin are enriched for association with age



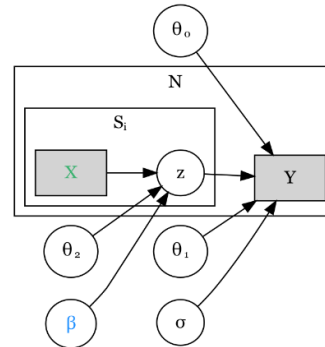
B Intronic and intergenic mutations are enriched for association with age



C Deleterious mutations are depleted of association with age



D The genomic and epigenomic mutation rate (GEM)



212

213 The Genetic Determinants of Mutation Burden

214

215 Next, as clonal phenomena have been shown to have germline genetic determinants, we  
 216 performed a GWAS with GEM as the phenotype in 51,399 carriers of diverse ancestry. We computed  
 217 summary statistics using SAIGE<sup>34</sup>. We detected six genome-wide significant loci, including *TERT*, *TCL1A*,

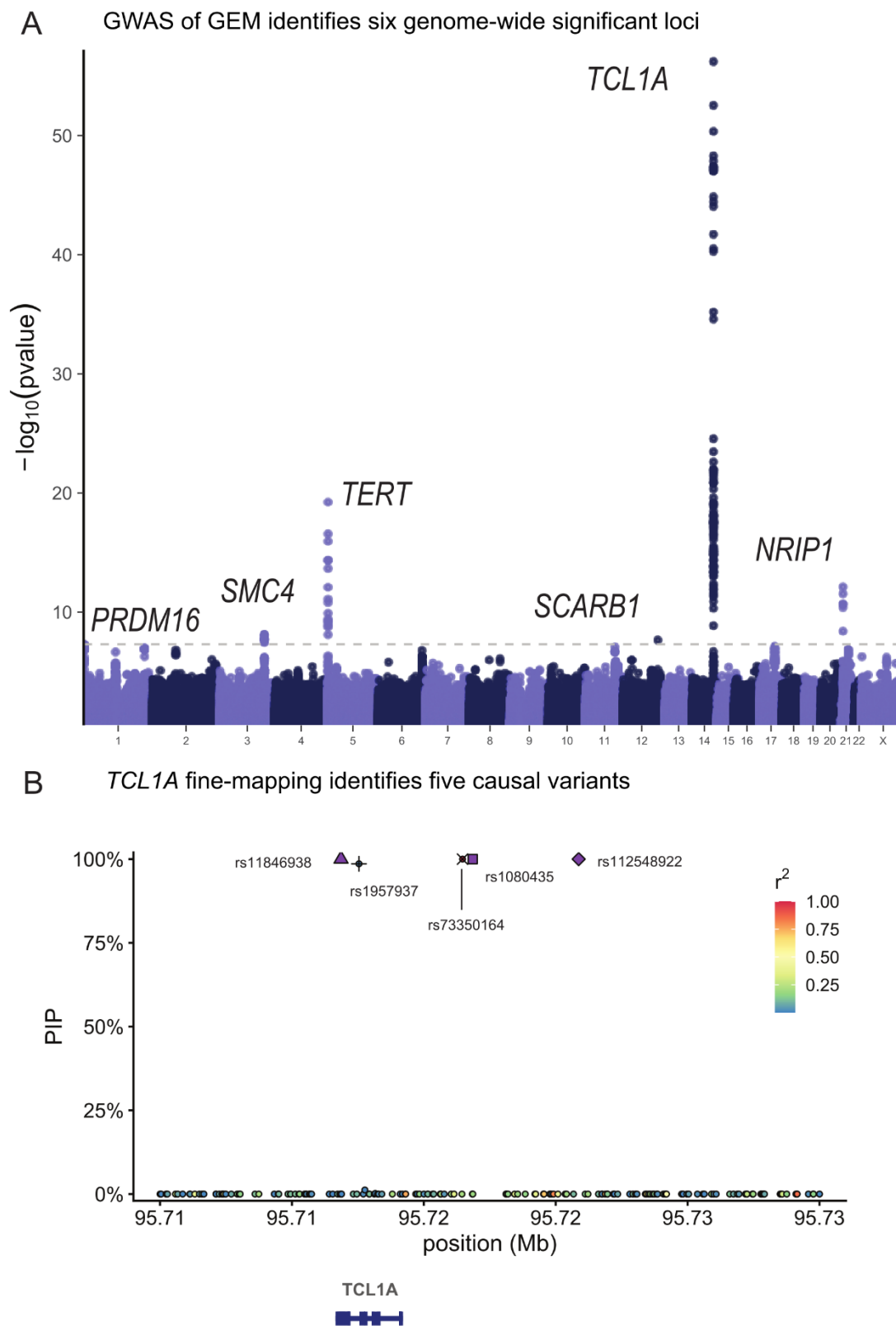
218 *TRIM59-SMC4-KPNA4, NRIP1, SCARB1, and PRMD16* (Figure 3A), and an overall  $h^2_{\text{SNP}}$  of 9.3%. We then  
219 performed a similar analysis based on the burden of mutations in quiescent chromatin and  
220 heterochromatin, which resulted in reduced power at *TCL1A* (GEM minimum pvalue of  $6 \times 10^{-57}$  vs  
221 mutation burden minimum pvalue of  $6 \times 10^{-50}$ ), demonstrating the value of the GEM over simply using  
222 mutation burden stratified by chromatin context (Extended Data Figure 2). *TERT, TCL1A, TRIM59-SMC4-*  
223 *KPNA4* have been previously associated with CHIP<sup>17,18,35</sup> and barcode-CH<sup>24</sup>, while *NRIP1* and *PRDM16*,  
224 *TERT*, and *TCL1A* have all been associated with mCAs<sup>21-23</sup> and barcode-CH<sup>24</sup>. *SCARB1* has not been  
225 previously reported with a related phenotype. To nominate causal SNPs and genes, we fine-mapped  
226 each locus using SuSIE<sup>36</sup> and cross-referenced the credible sets with the Open Targets V2G estimates<sup>37</sup>  
227 and cell-type specific enhancers-gene pairs from the activity by contact model<sup>38</sup>. These signals  
228 collectively highlight the convergence of germline variation influencing CHIP, mCAs, and clonal  
229 hematopoiesis without known drivers.

230 We then asked whether the association between these GWAS loci and GEM was mediated  
231 through an association with mCAs. Using a recently developed atlas of mCAs in TOPMed<sup>20</sup>, we excluded  
232 all mCAs or LOX carriers ( $n = 13,399$ ) and performed another GWAS as a sensitivity analysis. After  
233 filtering to variants that were genome-wide significant in either GWAS, we observed that the effect sizes  
234 were remarkably concordant ( $R^2 = 0.997$ , Extended Data Figure 3), suggesting that the GWAS of GEM is  
235 not merely mediated by the effect of mCAs/LOX.

236 The lead variant at *TCL1A* was rs2887399, which we previously discovered as among the primary  
237 mediators of clonal expansion in CHIP carriers<sup>39</sup>, and has been previously reported in GWAS of LOY<sup>21,22</sup>.  
238 Fine-mapping<sup>36</sup> revealed five credible sets, each with a single causal variant (rs11846938, rs112548922,  
239 rs1080435, rs1957937, and rs73350164, all PIP > .98). V2G<sup>37</sup> identified *TCL1A* as the mostly likely causal  
240 gene at each of the five causal SNPs. rs11846938 is 10bp from and in very high LD with rs2887399 (EUR  
241  $R^2 = 0.88$ , AFR  $R^2 = 0.92$ ), and both reside in the core promoter of *TCL1A*. Given high LD between  
242 rs2887399 and rs11846938, and the functional evidence we previously observed for the effect of  
243 rs2887399, we refer to this signal as “rs2887399/rs11846938.” rs112548922 is 9kb upstream of *TCL1A*,  
244 suggesting that cis-regulatory elements besides the core promoter are implicated in altered mutation  
245 burden. We previously described a CHIP-mutation specific mechanism, whereby mutations in  
246 *TET2/ASXL1/SF3B1* are associated with the aberrant chromatin opening of the *TCL1A* promoter in HSCs,  
247 but not *DNMT3A* mutations. This led us to hypothesize that rs2887399 is chromatin accessibility-QTL  
248 (caQTL) that is a *TET2/ASXL1/SF3B1* mutation (and possibly LOY) specific and thus the risk allele is more  
249 likely to lead to the aberrant activation of proto-oncogene *TCL1A* only when a mutation sufficient for  
250 chromatin modification at the *TCL1A* promoter has been acquired. Indeed, rs2887399 has been since  
251 reported as a caQTL<sup>40</sup> in lymphoblastoid cell lines (LCL). We previously showed that *TCL1A* expression is  
252 sufficient for promoting clonal expansion and altering stress response in HSCs<sup>39</sup> and is a key regulator of  
253 clonal expansion in CHIP clones. The discovery of additional casual variants at the *TCL1A* locus highlights  
254 the utility of applying GEM to samples unascertained for specific genetic lesions, providing increased  
255 power for genetic discovery and indicates that the context specific up-regulation of *TCL1A* is a broader  
256 phenomenon than previously appreciated, and likely occurs in HSCs without mCAs, LOY, or CHIP,  
257 possibly mediated through stochastic epigenetic phenomena that result in increased accessibility at cis-  
258 regulatory elements of *TCL1A*.



259 Figure 3: The genetic determinants of GEM. A, The GWAS of GEM. Summary statistics were estimated with SAIGE. B, Fine-  
260 mapping of the *TCL1A* locus. Note rs11846938 is 10bp from rs2887399. Fine-mapping was performed with SuSIE.



261

262 The lead variant at *NRIP1* is rs2229742, a common (MAF = 6%) missense variant (p.Arg448Gly)  
263 predicted to be deleterious by SIFT<sup>41</sup>. *NRIP1* has previously been discovered in the context of mCAs<sup>14</sup>  
264 and barcode-CH<sup>24</sup> but has not been discovered in CH-LPMneg. Fine-mapping of the locus identified one  
265 credible set containing three variants (rs2229742, rs2823020, rs2823025). The C allele of rs2229742 is  
266 associated with increased GEM burden (beta = 0.08 standard deviations, 72% PIP, pvalue =  $7.6 \times 10^{-13}$ ).  
267 *NRIP1* is a highly conserved (pLI<sup>42</sup> = 0.99) transcription co-regulator that is highly expressed in HSCs<sup>43</sup> and  
268 has been previously reported as a positive regulator of stemness in HSCs<sup>44</sup>. rs2229742 is strongly  
269 associated with multiple blood cell index GWAS<sup>45</sup> indicating that altered protein sequence of *NRIP1*  
270 results in altered HSC function. *NRIP1* ablation has previously been shown to extend lifespan in murine  
271 models<sup>46</sup>, indicating a role in aging related phenotypes, although this is reported to be likely mediated  
272 through the interaction between *NRIP1* and estrogen signaling rather than modulation of HSC function.  
273 As the C allele is associated with increased GEM, this suggests that the C allele may either increase  
274 function of the *NRIP1* product or increase the abundance of *NRIP1* through indirect mechanisms,  
275 perhaps through increased translation efficiency. To elucidate the consequences of altered amino acid  
276 sequence in *NRIP1*, we cross-referenced a recently released catalogue of *trans*-pQTLs from plasma<sup>47</sup>.  
277 rs2229742 is a *trans*-pQTL for both *SDC4* (beta = 0.09, pvalue =  $3.8 \times 10^{-13}$ ) and *PGLYPR2* (beta = 0.07,  
278 pvalue =  $8.7 \times 10^{-14}$ ). *SDC4* is a syndecan, which are cell-surface proteins that can interact with a broad  
279 range of ligands. The mouse-genome informatics resource<sup>48</sup> reported that *SDC4* ablation in mice led to  
280 several altered hematopoietic phenotypes. *PGLYPR2* is a peptidoglycan recognition protein that has  
281 been implicated in interferon regulation and innate immune response<sup>49</sup>.

282 Previous GWAS of clonal phenomena, including CHIP and MPNs have reported the *TRIM59*-  
283 *SMC4-KPNA4* locus, though none have conclusively identified the causal gene. Fine-mapping this locus  
284 identified a credible set containing 19 variants. The three variants with the highest PIP were rs11718121  
285 (PIP = 9.2%), rs1451760 (PIP = 8.8%), and rs6790951 (PIP = 8.2%). V2G estimated that *SMC4* was the  
286 mostly likely causal gene for each of these three variants, and alt-alleles at the three variants were  
287 associated with increased expression of *SMC4* in eQTLGen in whole blood<sup>50</sup> and increased expression of  
288 *SMC4* in lipopolysaccharide stimulated monocytes<sup>51</sup>. The interval spanned by the credible set contains a  
289 predicted *SMC4* enhancer in CD34+ cells by the ABC model. *SMC4* is a sub-unit of the condensin  
290 complex, which is involved in chromosome assembly and segregation during mitosis. Collectively, fine-  
291 mapping, V2G estimates, and the ABC model nominate *SMC4* as the mostly likely causal gene in the  
292 locus, highlighting the role of *SMC* related proteins in modulated mutation burden.

293 Fine-mapping of the *SCARB1* locus identified one credible set with a single SNP, rs11057853, a  
294 common (MAF = 46%) variant intronic to *SCARB1*. V2G estimated that *SCARB1* was the most likely causal  
295 gene for rs11057853, supported by both its presence within the *SCARB1* gene body and its role as an  
296 eQTL for *SCARB1* in blood<sup>50</sup>. The C allele was associated with reduced GEM (beta = -0.03 GEM standard  
297 deviations) and reduced expression of *SCARB1* (beta = -0.23, pvalue =  $5.5 \times 10^{-175}$ ), suggesting that  
298 *SCARB1* may be protective against CH-LPMneg. *SCARB1* is a receptor for HDL and rare variant burden  
299 tests of *SCARB1* in UK Biobank have identified several associations with lipid traits<sup>52</sup>. Previous reports  
300 have described a possible role for *SCARB1* in mediating the metabolic adaptation of long-term HSCs  
301 using murine models<sup>53</sup>.

302 We then asked whether rare variants are associated with GEM. We performed a genome-wide  
303 RVAS using STAAR<sup>54</sup>, including all missense and loss of function (LOF) variants within protein-coding  
304 genes. We identified 33 and 18 hits at pvalue thresholds of  $5 \times 10^{-6}$  and  $5 \times 10^{-7}$  (Supplementary table 3).  
305 The strongest hit was *CELF2* (pvalue =  $3.4 \times 10^{-28}$ ), where coding variants were associated with a higher  
306 GEM value among carriers (mean of 0.38, 95% CI: [0.36, 0.40]) than non-carriers (mean of 0.00, 95% CI:  
307 [-0.01, 0.01]). *CELF2* is a highly constrained (pLI<sup>42</sup> = 1.00) RNA binding protein and is highly expressed in

308 neutrophils<sup>37</sup>. A recent report described the role of *CELF2* as a suppressor of the AKT/PI3K signaling  
309 pathway<sup>55</sup> in lung carcinoma, which is presumably the signaling pathway mediating the effect of *TCL1A*.

310 We performed a non-coding RVAS using SCANG<sup>56</sup>, a dynamic window approach for identifying  
311 sets of SNPs that associate with phenotypes. We applied to SCANG to 100kb regions flanking 1,688  
312 genes that Open Targets<sup>37</sup> has identified as previously associated with cancer (Supplementary table 4).  
313 We identified 52 and 6 hits at pvalue thresholds of  $5 \times 10^{-6}$  and  $5 \times 10^{-7}$  (Supplementary table 5). We  
314 observed a locus on chr14 flanking the *MARK3* gene which was strongly associated with mutation  
315 burden ( $2.5 \times 10^{-8}$ ). To identify the likely causal variants within this window, we then performed a joint  
316 analysis of all rare variants that were included. This analysis highlighted three variants as significantly  
317 associated including rs190231639, a rare (MAF = 0.04%) variant intronic to *COA8*. To identify the likely  
318 causal gene in this gene dense locus, we cross-referenced the V2G results from Open Targets, which  
319 nominated *COA8*, *KLC1*, *XRCC3*, and *ZFYVE21* as equivalently likely causal genes. *XRCC3* is involved in the  
320 homologous recombination repair pathway of double-stranded DNA, though we are unable to  
321 conclusively identify it as the causal gene. Collectively, we identify diverse signals among the rare  
322 variants implicating DNA repair and post-translational modifications as key molecular processes  
323 contributing to GEM variation.

324 Mutation Burden Is Indirectly Regulated by the Size of the HSC Pool  
325

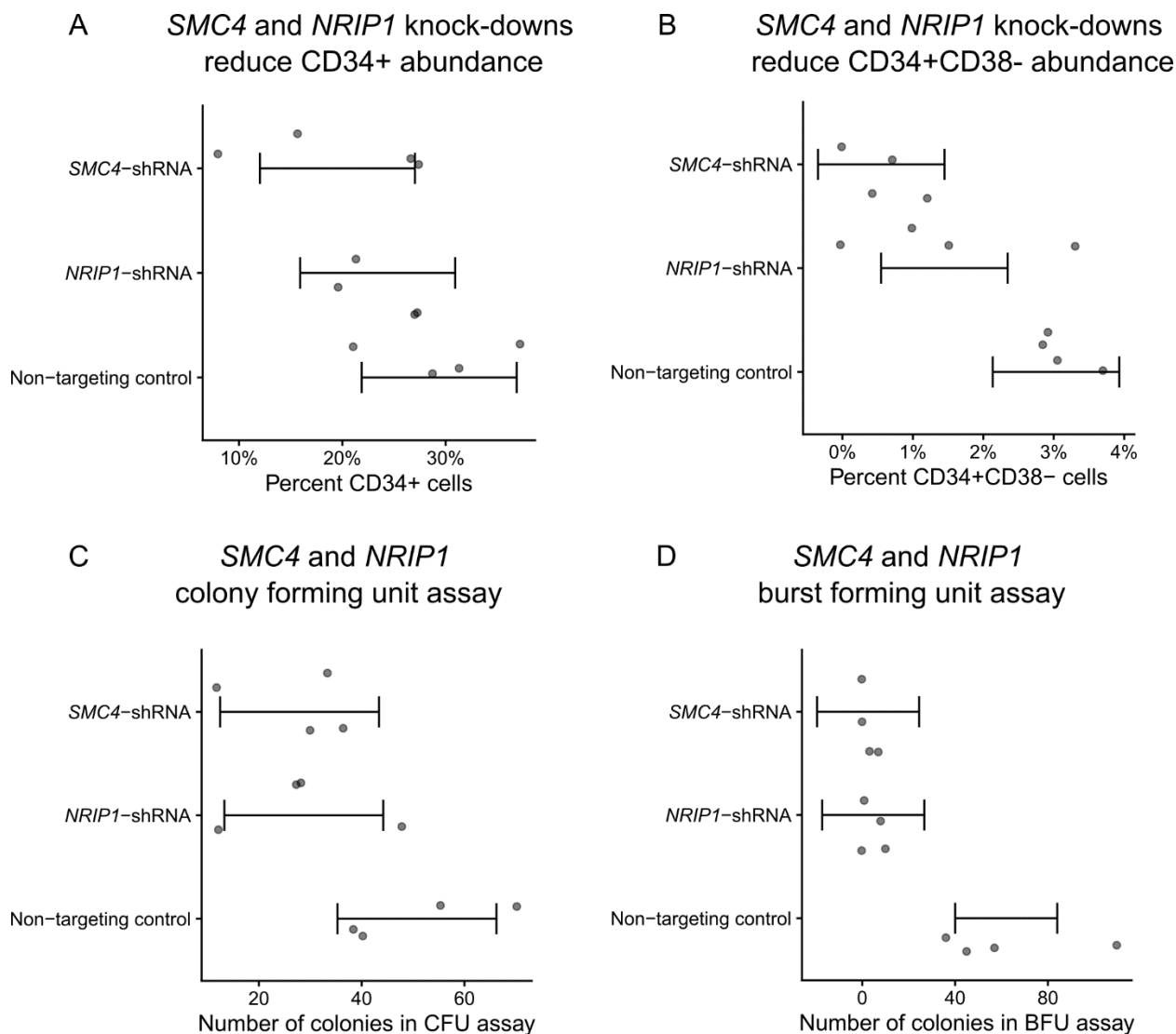
326 We then asked whether *SMC4* and *NRIP1* contributed to mutation burden by altering HSC self-  
327 renewal. To test this hypothesis, we separately knocked down *SMC4* and *NRIP1* in CD34+ bone-marrow  
328 derived human HSCs with shRNAs and performed a colony forming unit (CFU) assay. Relative to a non-  
329 targeting shRNA control, both *SMC4* and *NRIP1* knockdown cells were more likely to lose stemness in  
330 culture (1.6% and 2.5% fewer CD34+CD38- cells, pvalues =  $2.0 \times 10^{-3}$ ,  $1.7 \times 10^{-3}$ , Figure 4A-B) and the cells  
331 formed fewer burst-forming unity colonies (22.0 and 22.8 fewer colonies, pvalues =  $2.4 \times 10^{-3}$ ,  $1.9 \times 10^{-3}$ ,  
332 Figure 4C-D). These results are consistent with a role of both *SMC4* and *NRIP1* as positive regulators of  
333 CD34+ HSC self-renewal, and led us to propose the following mechanisms for their roles as indirect  
334 regulators of mutation burden: Either *SMC4/NRIP1* may regulate the fitness of HSC in response to  
335 noxious stimuli such as the infection inherent to the shRNA knockdown assay, or *SMC4/NRIP1* directly  
336 regulates the size of the HSC pool; a larger active HSC pool increases the likelihood that at least one HSC  
337 obtains a fitness advantage through either a genetic lesion or stochastic epigenetic phenomena, leading  
338 to a clonal expansion, which will increase the passenger count burden. Similar models have been  
339 previously been proposed in the context of myeloproliferative neoplasms<sup>27</sup>.

340 We then sought to explore the causal relationship between HSC pool size and GEM through  
341 simulation using a stochastic process that describes realistic HSC population (Methods). We simulated  
342 several HSCs which acquire passenger mutations at a constant rate per cell through a Poisson point  
343 process. We simulated the size of individual HSC clone populations using a Poisson birth-death process,  
344 where a single parameter  $s$  governs the relative likelihood of an HSC self-renewing into two identical  
345 HSCs as opposed to dividing into two differentiated cell types. At a rate of 1 driver mutation per 10,000  
346 HSCs per year, we simulated modest increases to  $s$  in each HSC to model modest increases in cell fitness  
347 that some clones may acquire. We stratified these simulations across varying initial sizes of the HSC  
348 pool, finding that larger pools were much more likely to contain at least one clone with a substantial  
349 increase in fitness (Extended Data Figure 4) and many more high-VAF passengers (Extended Data Figure  
350 5). Importantly, the burden of high VAF passengers increased as the number of increases to  $s$  increased.  
351 Taken together, this model provides a formal exposition for why GEM may be associated with both the

352 overall HSC pool and the likelihood of at least one clone obtaining a substantial increase in self-renewal  
353 capacity.

354

355 *Figure 4: The functional consequence of SMC4 and NRIP1 on HSCs. A, SMC4 and NRIP1 were knocked-down with shRNA and the*  
356 *proportion of CD34+ cells was quantified with FACS. Quantities were compared referent to a non-targeting control. B, proportion*  
357 *of CD34+CD38- was quantified with FACS. C, Number of colonies formed in a colony-forming unit (CFU) assay. D, Number of*  
358 *colonies in a burst-forming unit assay.*



359

360 Regional Mutation Burden and the Genetic Determinants of Sex Specific Mutation Pathways

361

362 We then sought to perform more granular analyses of mutation burden, examining  
363 heterogeneity by sex and position in the genome. As the importance of mutation burden may vary  
364 based on genome position, we estimated the mutation burden in 49 non-overlapping intervals  
365 approximately  $5 \times 10^7$  bases in length. We then estimated the association between age and mutation  
366 count stratified by interval, indicating heterogeneous associations across the intervals (Extended Data  
367 Figure 6). To characterize the underlying structure, we then performed PCA on these mutation counts.

368 We observed that PC1 was strongly associated with overall mutation burden and explained 60% of the  
369 variance, revealing a single general factor associated with mutation burden genome-wide. We observed  
370 that the loadings of PC2 were enriched for mutations appearing chromosome-X, suggesting a sex-  
371 specific mutation pathway. We observed that PC2 is significantly associated with genotype inferred sex  
372 ( $R^2 = 7.3\%$ ,  $pvalue < 2.2 \times 10^{-16}$ ), highlighting a sex-specific contribution to somatic variation in whole  
373 blood.

374 We then asked whether this sex-specific mutation factor had distinct germline determinants  
375 from GEM. We observed that a single locus on chromosome 1 near *SRGAP2C* was associated with PC2.  
376 The lead variant at this locus is rs61804016, a common variant 62kb away from the transcription start  
377 sites of *SRGAP2C* that has been previously reported as an eQTL in whole blood for *SRGAP2C*, *NBPF8*,  
378 *NBPF26*, *PFN1P2*, and *SRGAP2C*. However, in monocytes and T-cells, rs61804016 is only an eQTL for  
379 *SRGAP2C*<sup>57</sup>. The eQTL associations and proximity to the TSS of *SRGAP2C* suggest that *SRGAP2C* is the  
380 most likely causal gene in the locus. *SRGAP2C* is a GTPase activating protein that is expressed in  
381 hematopoietic progenitor cells<sup>43</sup>. We then cross-referenced phewas<sup>58</sup> results in UK Biobank<sup>59</sup> and  
382 FinnGen<sup>60</sup>. The C allele of rs61804016 is nominally associated with increased risk for breast cancer  
383 (odds-ratios of 1.08, 1.07,  $pvalues$  of  $1.2 \times 10^{-5}$ ,  $3.2 \times 10^{-6}$ ), further supporting the sex-specific nature of  
384 PC2. *SRGAP2C* is on chromosome 1, suggesting sex-specific regulation of an autosomal gene in the  
385 genesis of sex-specific mutation burden. No SNP near *SRGAP2C* was associated at genome-wide  
386 significance with the GEM phenotype. These analyses highlight the value of subtyping in mutation  
387 burden estimation by revealing sex-specific factors.

388 Given findings of sex-specific mutation pathways, we then performed sex-stratified GWAS of  
389 GEM (Figure 5), which revealed two sex specific signals. At the *SMC4* locus, rs11718121 was associated  
390 with increased mutation count in females ( $\beta = 0.043$ ,  $pvalue = 5.1 \times 10^{-9}$ ) but much more weakly  
391 associated in males ( $\beta = 0.019$ ,  $pvalue = 0.038$ ). At a locus not identified in the standard analysis, near  
392 *MSRA*, rs117344298 was associated with decreased mutation count in males ( $\beta = -0.17$ ,  $pvalue = 2.4$   
393  $\times 10^{-8}$ ), but unassociated in females ( $\beta = -0.013$ ,  $pvalue = 0.58$ ). Other signals were largely shared  
394 between males and females. Cross-referencing of sex-biased eQTLs reported in GTEx<sup>61</sup> found that *MSRA*  
395 had nominally significant sex-biased eQTLs in tibial nerve tissue and Brain cortex, but not in whole  
396 blood. *SMC4* did not have sex-biased eQTLs, which may be the result of limited power to detect such  
397 effects.

398

399

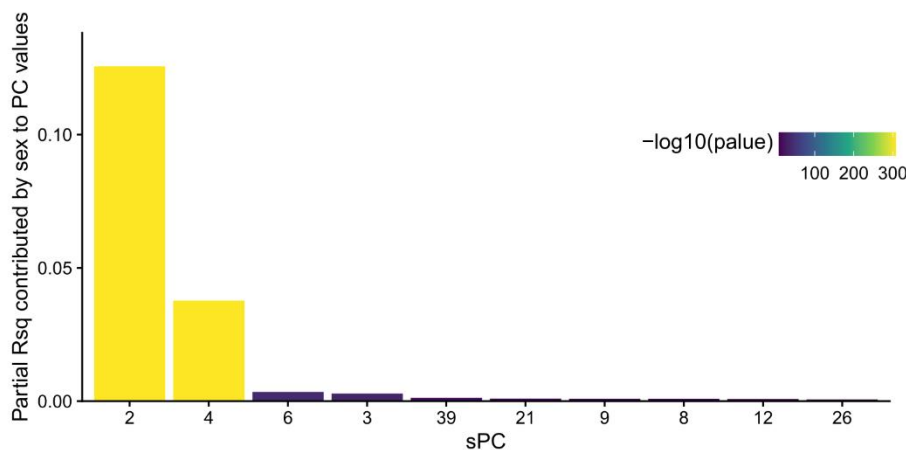
400

401

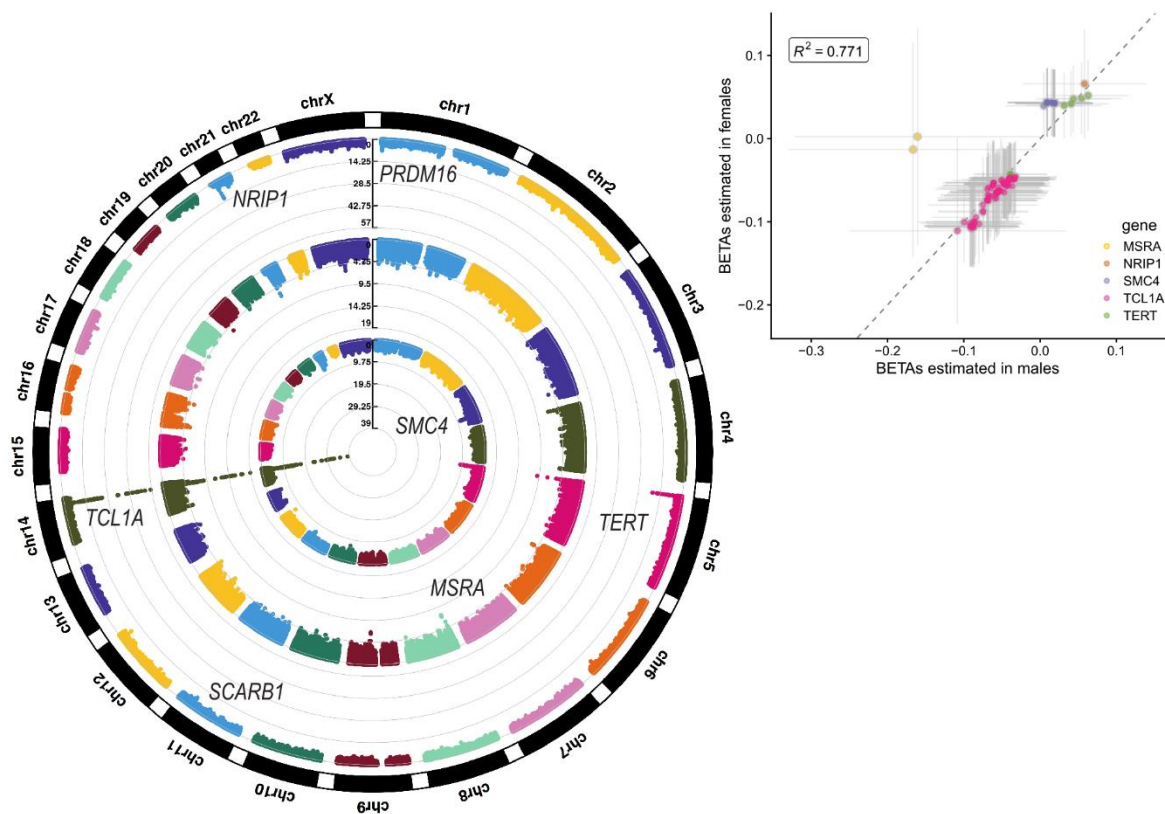
402

403 *Figure 5: The sex specific genetic determinants of mutation burden. A, Regressions were performed for each quantile-*  
 404 *transformed somatic principal component (sPC) on study and sex as covariates. The partial variance explained by sex is*  
 405 *displayed on the y-axis. B, Circular Manhattan plot. The outer-most ring is the GWAS of GEM on all individuals, the middle ring is*  
 406 *the GWAS of GEM on males, and the inner-most ring is the GWAS of GEM in females. Inset, a scatter plot of the two sex-specific*  
 407 *GWAS plotting all SNPs with  $p$  values  $< 1 \times 10^{-8}$  in either GWAS. Asymptotic confidence intervals are plotted with a width*  
 408 *corresponding to genome-wide significance.*

**A Somatic PC2 is associated with genotype inferred sex**



**B Sex stratified GWAS identifies sex-specific CH genes**



## 410 The Transcriptomic Correlates of High Mutation Burden

411

412 Next, we asked whether GEM associated with altered gene expression across five tissue types  
413 available in TOPMed, including whole blood, PBMC, monocytes, T cells, and nasal epithelial cells. We  
414 performed a search for GEM-gene associations by regressing the inverse normal transformed expression  
415 values of each gene ( $n = 17,741$ ) on the inverse normalized GEM estimates, including age, genotype  
416 inferred sex, 15 genotype PCs, 20 expression PCs, and cohort indicators as covariates. To increase  
417 power, we then used mashr<sup>62</sup> to apply shrinkage across the 88,705 GEM-gene associations. We  
418 identified 404 GEM-gene associations at a local false sign rate (lfsr)<sup>63</sup>  $< 0.05$  (Supplementary Table 6).  
419 Within whole blood, we observed the up-regulation of *RUFY4* with increased GEM. *RUFY4* is highly  
420 expressed in dendritic cells<sup>37</sup> and is involved in response to the anti-inflammatory cytokine IL-4<sup>49</sup> (Figure  
421 6A). We also observed the down-regulation of *LGSN*, which although annotated for its role in  
422 differentiation cells in the lens, is highly expressed in HSCs and was recently reported as a candidate  
423 causal gene in asthma<sup>64</sup>.

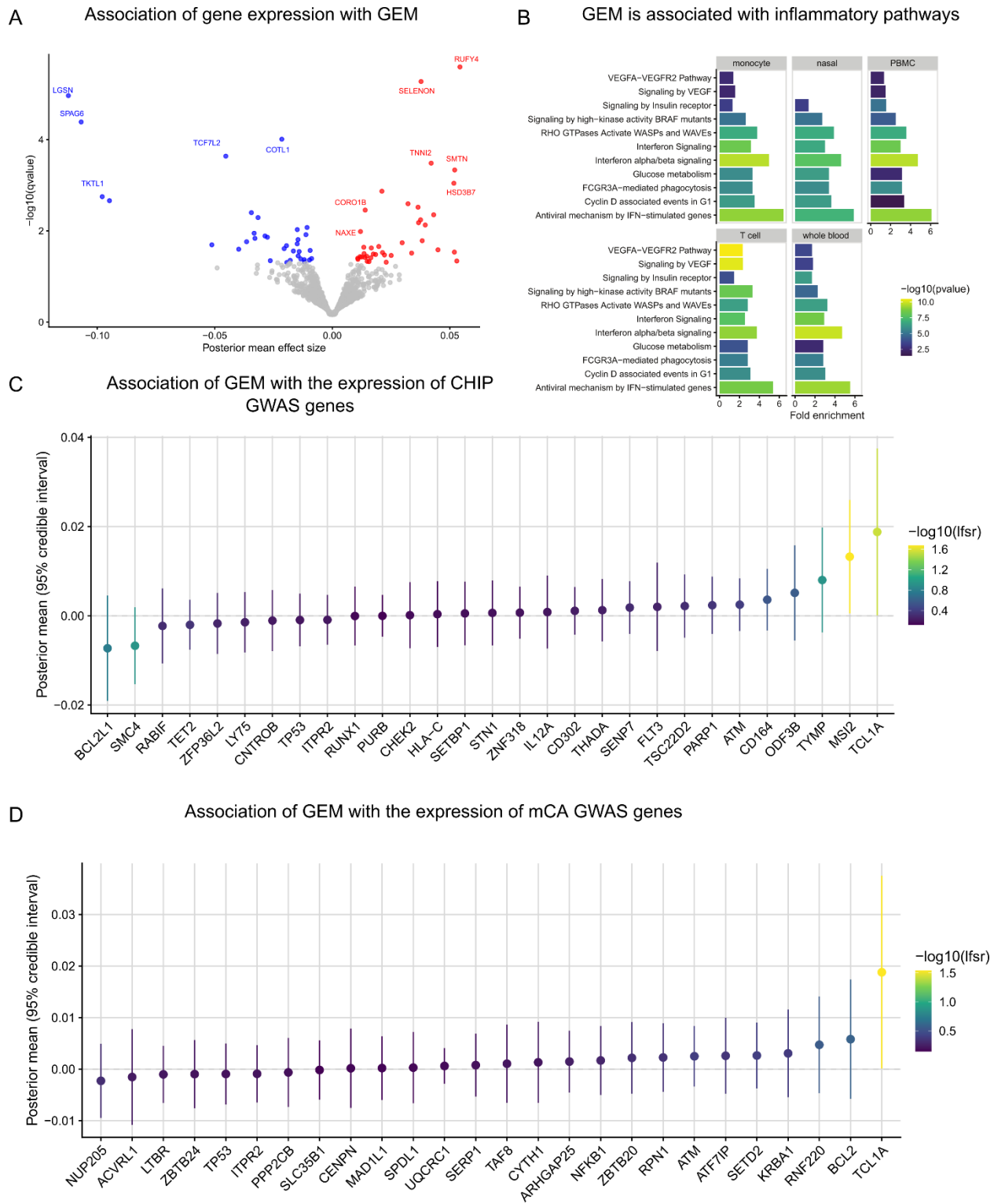
424 To characterize the gene programs associated with GEM, we performed pathway enrichment  
425 analyses using KEGG<sup>65</sup> as a reference. We observed a striking down-regulation of genes involved in  
426 interferon signaling (Fig. 6B). Interferon-alpha is a cytokine with anti-proliferative properties that was  
427 previously considered as a candidate therapeutic for AML<sup>66</sup>. Interferon-alpha is thought to reduce clonal  
428 expansion through direct and indirect mechanisms, including inducing apoptosis and activating the  
429 adaptive immune system. We then performed a tissue specific pathway analysis among the 404 GEM-  
430 genes, which similarly identified interferon alpha/beta signaling as greater than 3x fold enriched in each  
431 of the five tissue types. We also identified enrichment of the VEGFA-VEGFR2 pathway that is largely  
432 specific to T cells (Fig. 6B). Deletion of VEGFA in CD8+ T cells has previously been shown to reduce  
433 effector function<sup>67,68</sup>. Collectively, these results highlight the importance of anti-proliferative cytokines  
434 to inhibiting clonal expansion and suggest that transcriptomic responses to mutation burden in  
435 disparate tissues (nasal epithelial and blood samples) are more similar than anticipated.

436 Next, we asked whether specific loci that either define CHIP, or have been discovered in GWAS  
437 of CHIP or mCAs (Supplementary Tables 7-8), implicate genes whose expression levels are also  
438 associated with GEM. Among loci identified in either CHIP or mCA GWAS, we observed that expression  
439 of *TCL1A* and *MSI2* are positively associated with GEM (Fig. 6C-D). The association between *TCL1A* and  
440 GEM is consistent with *TCL1A* expression as a key mechanism in modulating mutation burden and clonal  
441 expansion in whole blood. Among CHIP mutations (Supplementary Table 9), we observed that *IDH2*  
442 expression is negatively associated with GEM (Extended Data Figure 7), which corroborates its  
443 protective effects against clonal expansion. We then asked whether within blood, there were specific  
444 genes with heterogeneous effects. We found that although effect sizes generally were highly  
445 concordant, *TCL1A* had a much stronger association with GEM in T cells than in monocytes (Extended  
446 Data Figure 8), highlighting the need the tissue and cell-specific transcriptomic analyses when  
447 performing searches in blood for the transcriptomic correlates of mutation burden.

448

449

450 *Figure 6: The transcriptomic correlates of GEM. A, Association analyses were performed between GEM and gene expression in*  
 451 *whole blood, including age, sex, genotype PCs 1-5, and expression PCs 1-20 as covariates. B, Enrichment analyses were*  
 452 *performed using pathfindR and KEGG pathways as reference. C, Association statistics among CHIP GWAS genes. D, Association*  
 453 *statistics among mCA GWAS genes.*





455 The Clinical Correlates of GEM  
456

457 Because clonal hematopoiesis phenomena have been previously associated with cardiovascular  
458 disease<sup>4,6,15</sup>, we asked whether GEM associates with vascular and heart disease phenotypes in TOPMed.  
459 We first asked whether GEM was associated with coronary artery disease (CAD). We restricted our  
460 analyses to those NHLBI TOPMed cohorts with harmonized longitudinal assessment of CAD events; this  
461 enabled separate analyses of the association between GEM and incident CAD phenotypes (i.e., CAD  
462 events that occurred after the blood draw from which GEM was assessed) and the association between  
463 GEM and prevalent CAD (i.e., CAD events prior to the GEM blood draw). Within four NHLBI cohorts  
464 (WHI, FHS, CHS, COPDGene), we performed a Cox-proportional hazards regression analyses for incident  
465 CAD events after excluding individuals with prevalent CAD disease, including GEM, age at baseline,  
466 smoking history, body-mass index (BMI), sex, and germline genotype PCs as covariates. We observed  
467 that GEM was not associated with incident CAD events (meta-analysis hazard ratio: 1.00, 95% CI: [0.97,  
468 1.04], pvalue = 0.84, Figure 7A).

469 We performed similar analyses with incident ischemic stroke. After meta-analyzing results from  
470 three cohorts (WHI, CHS, ARIC), we observed that there was substantial heterogeneity ( $I^2 = 86\%$ ) in  
471 results, with a positive association observed in WHI (hazard ratio of 1.19, 95% CI: [1.12, 1.26], pvalue =  
472  $2.8 \times 10^{-8}$ ) and null or negative effects observed in ARIC and CHS (Extended Data Figure 9). Meta-analysis  
473 resulted in no association between incident stroke and GEM (hazard ratio of 1.03, 95% CI: [0.88, 1.22]),  
474 Figure 7A). Given the differences in the distribution of sex across the three cohorts, we then performed  
475 a female only analysis in CHS and ARIC, finding no evidence for a sex-specific effect after meta-analyzing  
476 with WHI (hazard ratio of 1.05, 95% CI: [0.87, 1.26], Extended Data Figure 10). We then asked whether  
477 GEM is associated with incident peripheral artery disease (PAD). After meta-analysis, we observed that  
478 GEM was associated with increased risk (hazard ratio 1.15, 95% CI: [1.05, 1.26], pvalue =  $3 \times 10^{-3}$ , Figure  
479 7A) for incident PAD events.

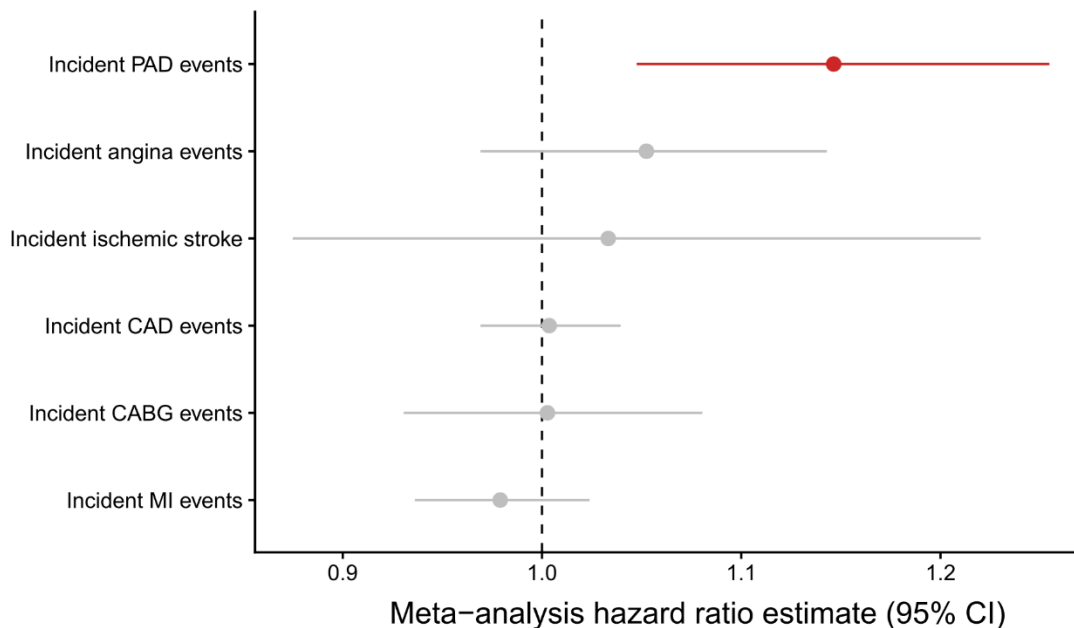
480 Because evidence within the clonal hematopoiesis literature (Heyde et al.<sup>69</sup>) suggests that prior  
481 CAD events are causal contributors to the up-regulation of HSC proliferation, we then asked whether  
482 prior CAD was associated with GEM. In a meta-analysis across five NHLBI cohorts (WHI, FHS, CHS, COPD,  
483 ARIC), we observed that prior CAD was suggestively associated with a modest increase in GEM (effect of  
484 prior CAD on standardized GEM: 0.12, 95% CI: [-0.03, 0.28], pvalue = 0.13, Extended Data Figure 11).

485 To reconcile these disparate phenotypic associations, we performed analyses between GEM and  
486 biomarkers, including complete blood cell counts (CBC) measurements and inflammation  
487 measurements. We observed that GEM was positively associated with increased white blood cell count  
488 after adjustment for age and smoking status at baseline (Fig. 7B). These results indicate that GEM is  
489 associated with altered hematopoiesis. In contrast, after meta-analysis, we observed no association with  
490 CRP (Figure 7B). This suggests that the association between GEM and PAD is not mediated through  
491 systemic inflammatory markers like CRP, but we note that markers of systemic inflammation like CRP  
492 may not be sufficiently sensitive to capture the association between inflammation and HSC activity  
493 within the bone-marrow microenvironment.

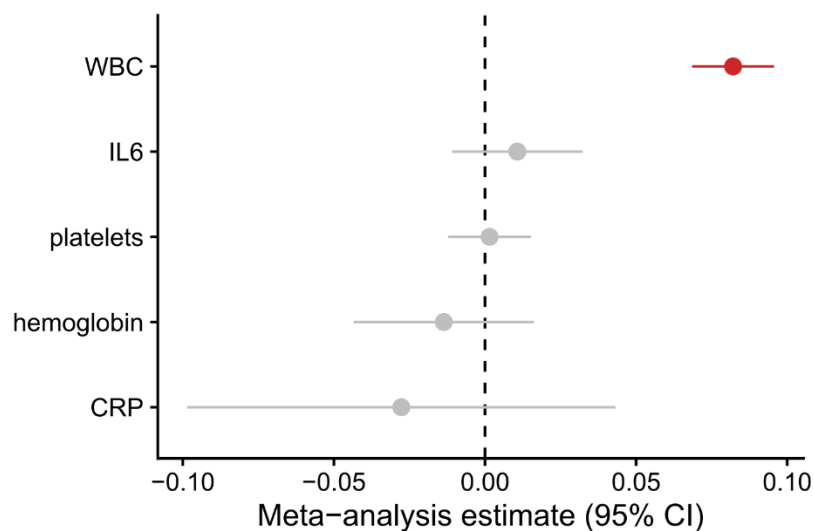
494

495 *Figure 7: The phenotype correlates of GEM. A, Cox proportional-hazard regressions were performed, regressing incident events*  
496 *on a spline of age, sex, smoking status, and germline PCs. Individuals with prevalent disease were excluded. CAD = coronary*  
497 *artery disease, PAD = peripheral artery disease, CABG = coronary artery bypass graft, MI = myocardial infarction. CAD events*  
498 *were defined as at least one of an MI, CABG, angina, or angioplasty during the follow-up period. A random effects meta-analysis*  
499 *was performed. GEM was inverse normal transformed. Sex was excluded from the WHI regression, and smoking was excluded*  
500 *from the COPD regression. B, A linear regression of the inverse normal transformed biomarker, including a spline of age, sex,*  
501 *smoking status, and germline PCs as covariates. GEM was inverse normal transformed. Sex was excluded from the WHI*  
502 *regression, and smoking was excluding from the COPD regression.*

### A GEM is associated with increased risk for incident peripheral artery disease



### B GEM is associated with increased white blood cell counts



503

504

505 Conclusion

506 Using 51,399 diverse TOPMed whole genomes, we derived a semi-supervised model of mutation  
507 rate, GEM, that increases power for discovery of germline determinants of mutation burden by better  
508 distinguishing somatic mutations from sequencing artifacts. Using GEM to identify the germline  
509 determinants of CH-LPMneg, we observed the convergence of common variant loci influencing multiple  
510 types of clonal phenomena, including CHIP and mCAs, demonstrating that genetic predisposition to  
511 mutation burden is shared across several different clonal contexts. We observed that altered protein  
512 sequence of *NRIP1* was strongly associated with increased mutation rate, which along with its  
513 documented role in GWAS of blood cell indices, collectively implicate *NRIP1* as an important regulator of  
514 aberrant hematopoiesis. We also observed that *TCL1A*, which we previously identified as a critical  
515 moderator of clonal expansion in CHIP carriers, is also associated with mutation burden in samples  
516 ascertained for not having CHIP, which may reflect its contribution towards clonal expansion in other  
517 kinds of clonal phenomena. Using a sex specific mutation pathway revealed by PCA analysis, we  
518 identified a breast-cancer locus, *SRGAP2C*, that associates with GEM.

519 We anticipate that our approach, which identifies distinct underlying mutation pathways  
520 through PCA analysis can be extended to identify other factors that may contribute to specific mutation  
521 pathways. Overall, our approach identifies several loci that have been previously discovered in other  
522 analyses of clonal hematopoiesis phenomena, suggesting that studying clonal hematopoiesis without  
523 known drivers represents an under-appreciated model for discovering the germline determinants of  
524 mosaicism in blood. Importantly, analysis of mutation burden without a known CHIP genetic lesion  
525 greatly expands the sample size available to perform these analyses; our analysis here is an order of  
526 magnitude larger than the number of CHIP carriers discovered in our previous analyses of TOPMed  
527 analyses<sup>13,17</sup>.

528 We performed the first multi-tissue analysis of the transcriptomic consequences of mutation  
529 burden in whole blood. We observed the striking down-regulation of interferon signaling across five  
530 tissues, including four from blood and one from epithelial tissue. Collectively, this analysis highlighted  
531 the need for additional characterization of the *in-vivo* transcriptomic consequences of anti-proliferative  
532 cytokines on HSC growth. Interferon-alpha, among other cytokines with similar effects, may represent  
533 candidates for therapeutic intervention.

534 We observed that mutation burden in whole blood was associated with altered blood cell  
535 indices and increased risk for peripheral artery disease. However, GEM was not associated with incident  
536 CAD events. This is consistent with the observation that the association between CH and CAD is highly  
537 heterogenous across different forms of clonal phenomena. Within CHIP, the largest phenotype analysis  
538 to date<sup>35</sup> reported an association (1.31 hazard ratio) with *TET2* CHIP but not *DNMT3A* CHIP. Several  
539 analyses within smaller cohorts have also reported associations between CHIP and CAD phenotypes<sup>4,6,70</sup>.  
540 CH mediated through mCAs have no reported association with CAD<sup>9</sup>, while both positive and negative  
541 reports exists regarding LOY and CAD related phenotypes<sup>15,23</sup>. A recent analysis that examined barcode-  
542 CH<sup>24</sup>, which includes several different forms of CH, reported no association between barcode-CH and  
543 CAD, while finding an association between barcode-CH and PAD, concordant with our results. These  
544 observations reflect the need to examine the associations between CH and CAD stratified by the  
545 particular genetic lesion(s), size of clone, and potentially the rate at which a clone is expanding. Indeed,  
546 recent reports<sup>52,71</sup> on the plasma proteomic associates of CH have found substantial heterogeneity  
547 across different CHIP mutations, highlighting the substantial heterogeneity observed at both  
548 epidemiologic and molecular levels.

549 Our approach is not without limitations. The precise estimation of mutations from whole blood  
550 remains challenging. Although this approach has been shown to be promising in large cohorts in the  
551 context of epidemiologic association analyses, more sensitive sequencing assays are needed for clinical  
552 application. Additionally, the genesis of several mutations remains unclear. Although clonal expansion  
553 without known drivers clearly occurs<sup>12</sup>, elucidating the underlying mechanism remains an open  
554 question.

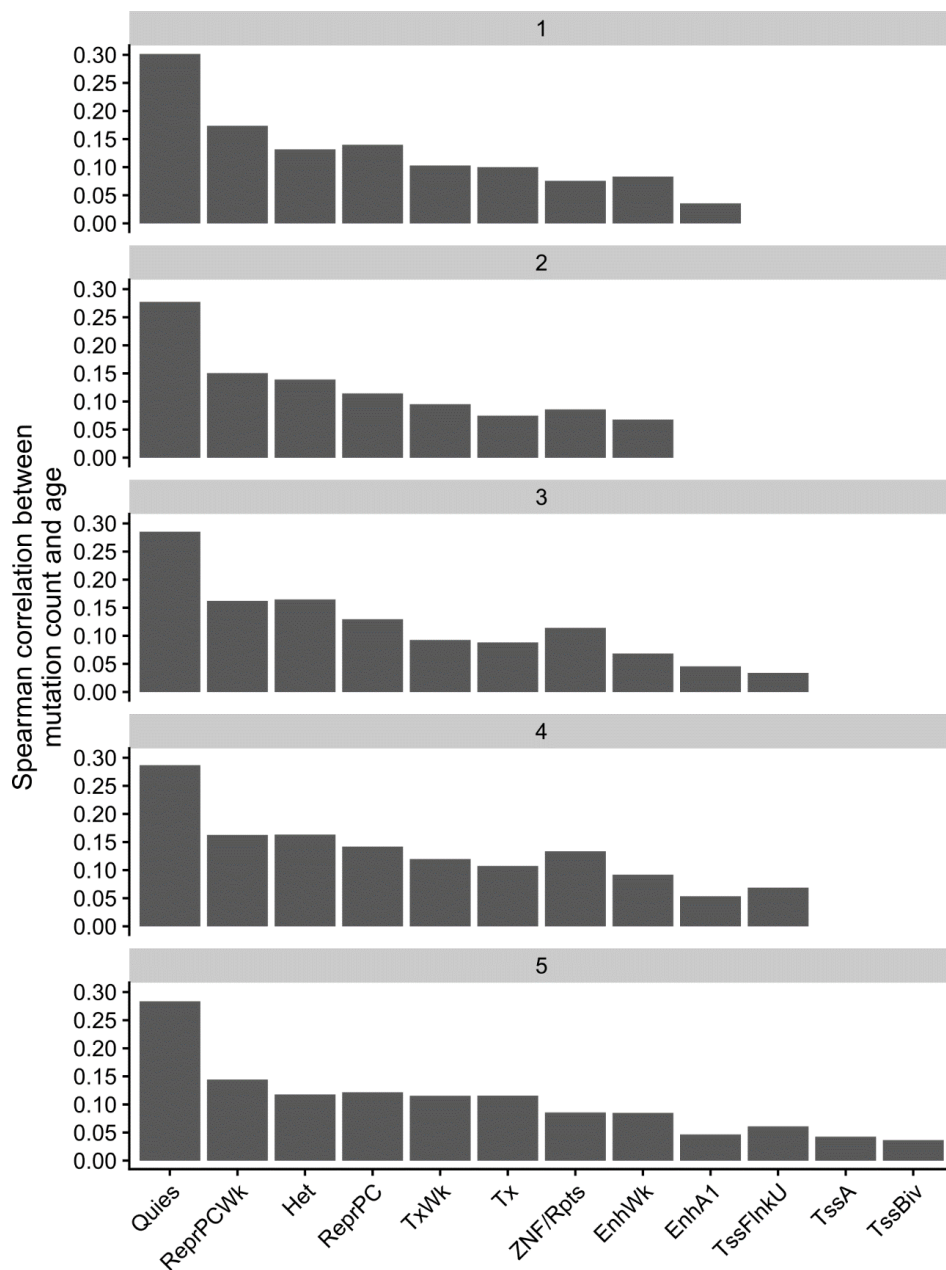
555 Overall, we develop a novel estimator of mutation burden that is not specific to CHIP carriers.  
556 We find that measuring mutation burden, even in individuals without known genetic lesions, is  
557 informative for aging related phenotypes. In contrast to surveillance for CHIP, which is relatively rare in  
558 individuals less than 80 years old, GEM can be used to monitor mutation burden in a larger proportion  
559 of adults. We anticipate that our approach will prove useful in non-blood tissues for the discovery of the  
560 germline basis of mutagenesis and will facilitate epidemiologic association analyses, ultimately  
561 elucidating the genesis and consequences of mutation burden.

562

563 Extended Data Figures

564 *Extended Data Figure 1: Spearman correlation between mutation burden and chronological age was*  
565 *calculated for each of the strata defined by chromHMM 15 state model in CD34+ cells and CADD derived*  
566 *quintiles. A CADD score of 5 indicates a score within the top 20% most deleterious variants.*

567



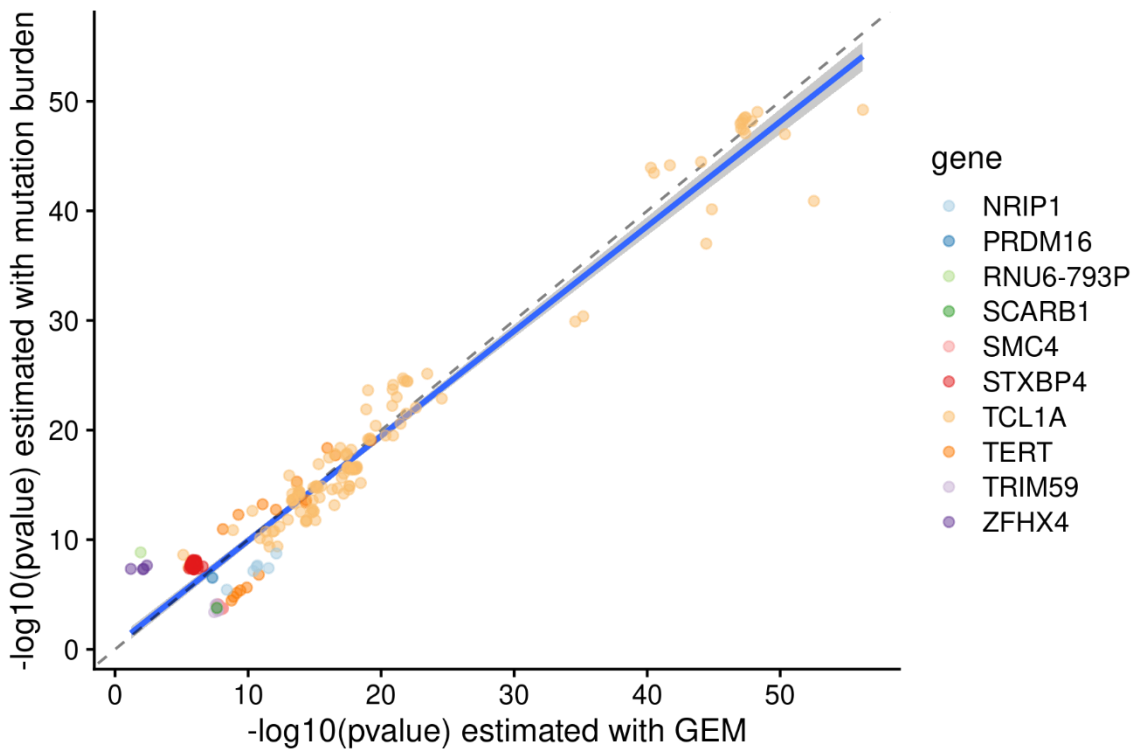
568

569

570

571

572 *Extended Data Figure 2: Scatter plot comparing the  $-\log_{10}$  pvalues from GWAS where the phenotype was*  
573 *either GEM (x-axis) or the burden of mutations falling in either heterochromatin or quiescent chromatin*  
574 *in CD34+ cells. Genes are colored by the likely causal gene, which was manually curated. Variants shown*  
575 *have pvalue  $< 5 \times 10^{-8}$  in at least one of the two GWAS.*



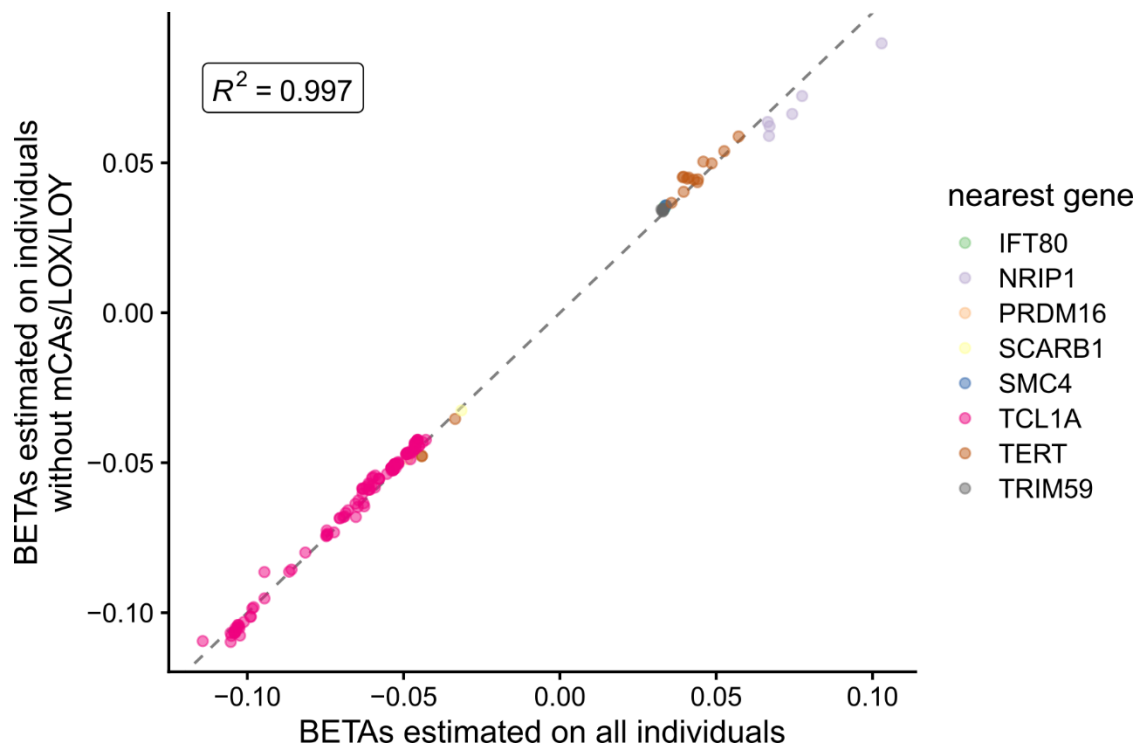
576

577 *Extended Data Figure 3: Scatter plot comparing the beta values from GWAS where the phenotype was*  
578 *either GEM on all individuals (x-axis, n = 51,399) or GEM on individuals that did not have an mCA (n =*  
579 *38,000). Genes are colored by the likely causal gene, which was manually curated. Variants shown have*  
580 *pvalue < 5 x 10<sup>-8</sup> in at least one of the two GWAS.*

581

582

583



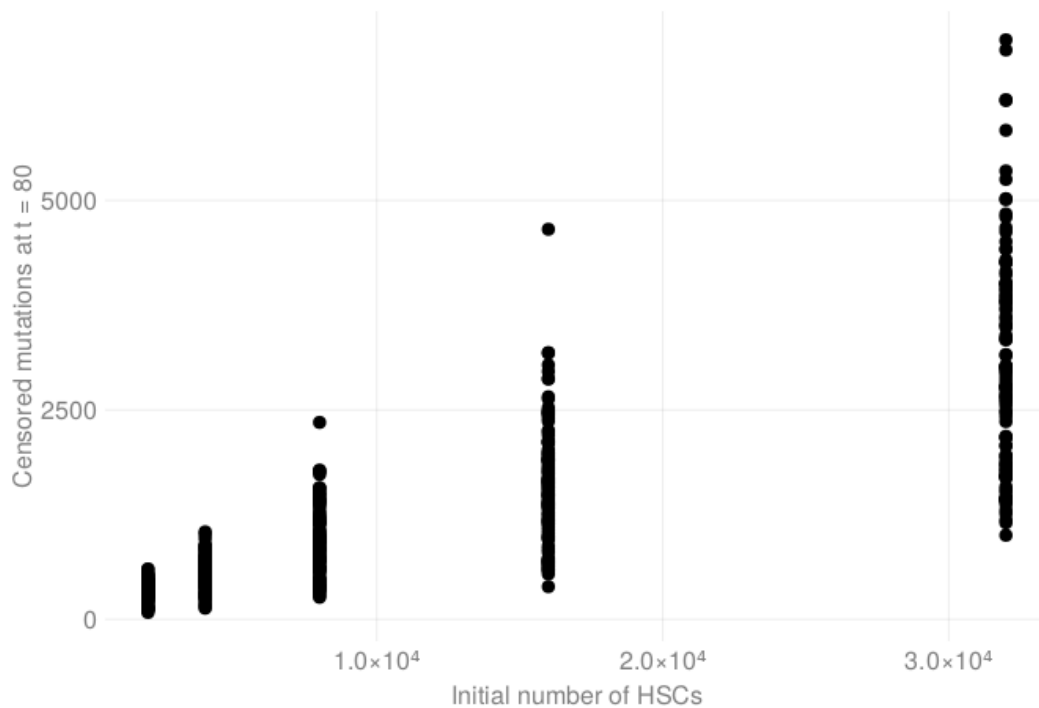
584

585

586

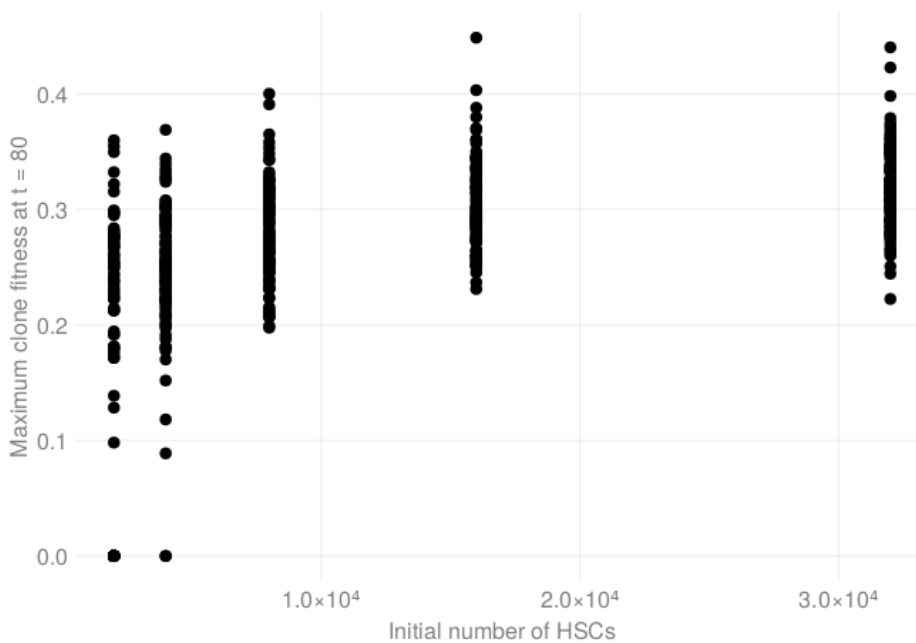
587 *Extended Data Figure 4: HSC stochastic process simulation, showing that the number of active HSCs has*  
588 *a large effect on the number of high-VAF mutations at the end of the simulation*

589



590

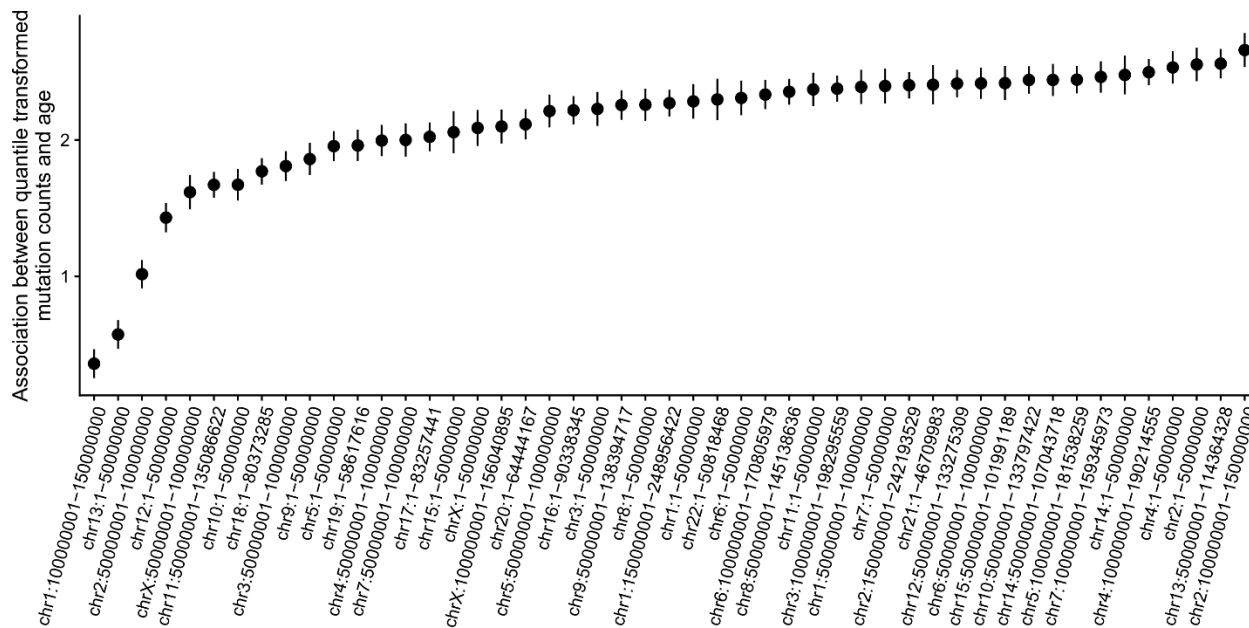
591 *Extended Data Figure 5: HSC stochastic process simulation, showing that the number of active HSCs has*  
592 *a large effect on likelihood of obtaining at clone with high fitness*



593



594 *Extended Data Figure 6: Linear regressions were performed between the inverse normal transformed*  
595 *mutation burden in each genomic bin with chronological age on the y-axis. Each regression include a*  
596 *study indicator as a covariate.*



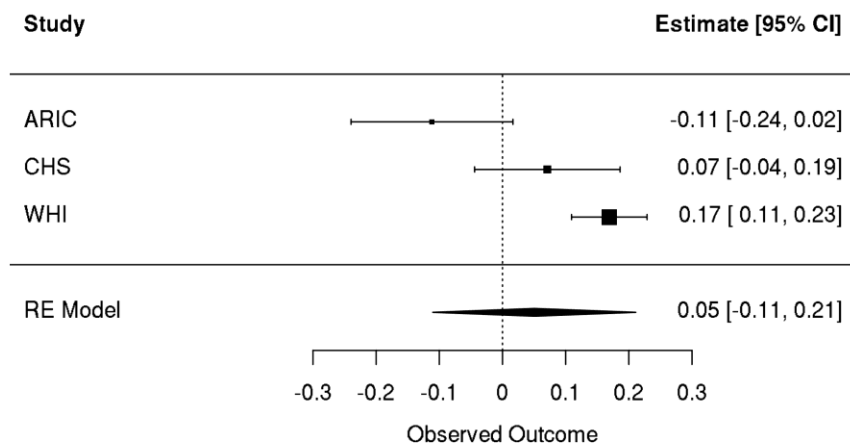
597

598

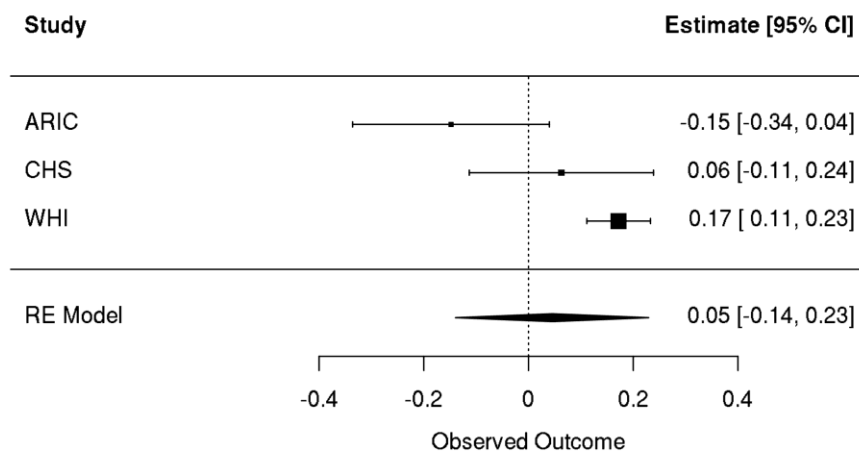
599



610 *Extended Data Figure 9: Meta-analyses of Cox proportional hazards regression with time to ischemic*  
611 *stroke as the outcome. A spline of age, sex, smoking status, and germline PCs were included as*  
612 *covariates. Individuals with prevalent disease were excluded.*

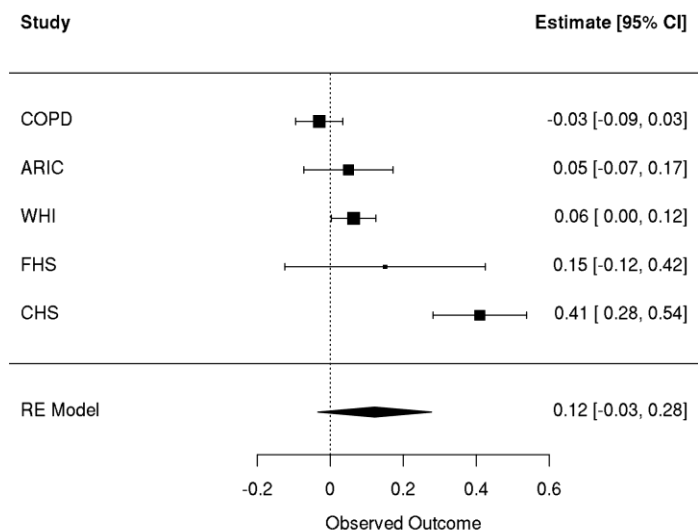


613  
614 *Extended Data Figure 10: Female only meta-analyses of Cox proportional hazards regression with time to*  
615 *ischemic stroke as the outcome. A spline of age, sex, smoking status, and germline PCs were included as*  
616 *covariates. Individuals with prevalent disease were excluded.*



617

618 *Extended Data Figure 11: Meta-analyses of linear regressions with inverse normal transformed GEM as*  
619 *the outcome and an indicator for prevalent coronary artery disease events that occurred prior to the*  
620 *blood draw that GEM uses as the covariate of interest. A spline of age, sex, smoking status, and germline*  
621 *PCs were included as covariates.*



622

623

## 624 Methods

625 Germline and somatic variant calling

626 TOPMed germline variant was performed as previously described<sup>25</sup>. Briefly, TOPMed BAM files were  
627 harmonized through the functionally equivalent pipeline<sup>72</sup>. Joint calling of germline SNPs and indels was  
628 performed with the Got-Cloud pipeline<sup>73</sup>. Samples were aligned to GRCh38. TOPMed germline SNP and  
629 indel freeze 10 was used in this analysis.

630 Putative somatic variants were first called with GATK Mutect2<sup>74</sup> in “tumor-only” mode with largely  
631 default settings. A “panel of normals” was included to exclude sequencing artifacts. Variant calling was  
632 performed on Google Cloud using Cromwell<sup>75</sup>. Only bi-allelic variants that passed Mutect2 filters were  
633 included in downstream analyses. CHIP calling was performed as previously described<sup>13,17</sup>; briefly, the  
634 Mutect2 output was cross-referencing with a list of predominately loss-of-function and missense  
635 mutations in a curated set of genes<sup>4,76</sup>.

636 We first identified somatic mutations that occurred once across all individuals, as singleton passenger  
637 mutations have a stronger association with chronological age than non-CHIP recurrent somatic  
638 mutations<sup>77</sup>. On mutations on the X-chromosome, we halved the variant allele-fraction for all mutations.  
639 We then excluded several mutations based on the following filters:

- 640 1. All mutations with a depth less than 25x or greater than 100x
- 641 2. All mutations falling within low complexity sequence regions
- 642 3. All mutations in segmental duplications
- 643 4. All mutations falling within genomic regions with germline CNVs with at least 10% minor-allele  
644 frequency. Germline CNVs from the TOPMed germline structural variant call-set<sup>78</sup> were used in  
645 this filter.
- 646 5. All mutations falling within the contigs with sequence that differed between hg19 and hg38, as  
647 defined by the “Hg19 diff” track in the UCSC genome table browser.
- 648 6. Any germline variant in TOPMed germline SNP and index freeze 10 (derived from 184,878 WGS)  
649 with a minor allele count of at least 10 and a variant allele fraction between .26 and .74
- 650 7. Any mutation with fewer than 2 alt reads or greater than 6 alt reads. At 38x, this corresponds to  
651 a VAF interval of 5%-16%

652 Annotation of somatic mutations

653 Singleton mutations after the above filters were first annotated with the variant effect predictor<sup>79</sup> (VEP)  
654 including the “—flag\_pick”, “—check\_existing”, “—canonical”, and “—flag\_pick” flags. A CADD<sup>33</sup> plugin  
655 was also included. CD34+ chromatin annotations were downloaded for sample BSS00233 from Roadmap  
656 epigenomics<sup>29</sup>. Mutations were also annotated with the correspond mutation type (e.g., C->T, G-T, etc.).

657 The genomic and epigenomic mutation rate (GEM)

658 GEM is a Bayesian graphical model with the following form. In the outlier layer, standardized  
659 chronological age (standardized with  $(\text{age} - 60) / 10$ ) is the outcome variable, denoted as  $Y_i$ . We note  
660 that GEM is a “weakly-supervised” model in the sense that while individual mutations are unlabeled, the  
661 entire training process is “supervised” by chronological age.  $Y_i$  conditional on the number of true  
662 mutations within an individual is assumed to follow a gaussian distribution. Each individual  $i$  has a  
663 candidate set of mutations  $S_i$  which were identified by the above filtering processes. Instead of using  
664 this raw count, we instead replace the count with the expectation of a Bernoulli random variable

665  $Z_{i,j}$  which denotes whether the  $j$ th mutation in the  $i$ th individual is a “true mutation” (i.e., takes a value  
666 of 1.0) or is an artifact (i.e., takes a value of 0). We include a non-linear transformation  $g$  to the sum  
667 over the true mutation burden. In practice,  $g(x) = \log_2(x)$  worked well.

$$668 \quad Y_i \sim N \left( \theta_0 + \theta_1 * g \left( \sum_{j \in S_i} Z_j \right), \sigma \right)$$

669 The expectation of this random variable is specified through an inverse-logit transformation, i.e.,

$$670 \quad Z_{j,l} \sim \text{Bernoulli} \left( \text{sigmoid}(\theta_2 + \mathbf{X}_j \boldsymbol{\beta}) \right)$$

$$671 \quad E(Z_{i,j}) = \text{sigmoid}(\theta_2 + \mathbf{X}_j \boldsymbol{\beta})$$

672 Where  $\mathbf{X}_j$  represents a length  $p$  vector of annotations for the  $j$ th mutation and  $\boldsymbol{\beta}$  is a length  $p$  random  
673 vector of weights.  $\theta_2$  is included as a bias or intercept term.

674 The above assumptions express the likelihood of GEM. The prior of GEM is specified as follows:

$$675 \quad \boldsymbol{\beta} \sim N_p(0, I)$$

$$676 \quad \sigma \sim \text{LogNormal}(0.0, 1.0)$$

$$677 \quad \theta_1 \sim \text{LogNormal}(\epsilon, 1.0)$$

678 Where  $\epsilon$  is in practice set to  $5 \times 10^{-3}$ . Inference was performed by optimizing the maximum a-posteriori  
679 objective using an ADAM optimizer. GEM is implemented in the torch package in R.

680 Within the matrix of mutation annotations  $\mathbf{X}$ , we include the following annotations:

- 681 1. VEP annotated variant impact
- 682 2. VEP “somatic” annotation
- 683 3. CADD\_PHRED score
- 684 4. The mutation type
- 685 5. The variant allele fraction
- 686 6. The chromatin state prediction

687 GEM was trained on 2,000 randomly sampled individuals with 186,277 total candidate mutations among  
688 them.

689 Genome-wide association studies with GEM

690 In the context of genome-wide association studies (GWAS), the phenotype was defined as the expected  
691 burden of “true” mutations, i.e.,  $\sum_{j \in S_i} E(Z_j)$  within the  $i$ th individual. This phenotype was inverse  
692 normal transformed. GWAS summary statistics were estimated with SAIGE on all germline variants  
693 where the minor allele count was at least 400 (i.e.,  $\text{MAF} \geq 0.4\%$ ) among the analyzed samples. Germline  
694 principal components 1-10, somatic principal components 3-4, genotype inferred sex, a cohort indicator,  
695 chronological age, average sequencing depth per sample, and the residual between the raw and  
696 estimated true mutation burden were included as covariates. Somatic principal components 1-2 were  
697 excluded as they are strongly associated with total mutation burden and sex respectively. Somatic

698 mutations that had previously been identified as recurrent<sup>77</sup> were excluded from the summary statistics.  
699 All germline variants with a milk-SVM threshold below -0.30 or an individual specific Hardy-Weinberg  
700 equilibrium -log<sub>10</sub> pvalue above 5.0 were excluded.

701 Rare-variant association studies with GEM

702 Rare-variant association studies (RVAS) were performed with the same GEM derived phenotype as the  
703 GWAS. We performed a non-coding RVAS by examining rare-variants within 100kb of cancer associated  
704 genes as defined by Open Targets<sup>37</sup> using SCANG<sup>56</sup>, which performs a scanning procedure for genome  
705 regions that contain a set of rare variants that associate with the phenotype. We similarly performed a  
706 genome-wide coding variant RVAS using STAAR, including any rare-variant annotated as having a  
707 “MODERATE” or “HIGH” impact on amino acid sequence by VEP.

708 Simulation of mutation burden

709 We assume an HSC can fall into one of three states:

- 710 1. HSCs can divide into two HSCs (“self-renewal”)
- 711 2. HSCs can divide into two differentiated cells
- 712 3. HSCs can divide into one HSC and one differentiated cell

713 For the purposes of simulating a stochastic process of HSC population size, we treat state 3 as irrelevant  
714 because it does not affect the total number of self-renewing HSCs.

715 We define the HSC clone birth rate as:  $\lambda_i(t) \sim \text{Poisson}(\omega * X_i(t) * (1 + s_i(t)) * dt)$  and the HSC  
716 clone death rate as  $\psi_i(t) \sim \text{Poisson}(\omega * X_i(t) * (1 - s_i(t)) * dt)$ .  $\omega$  is a parameter that controls the  
717 rate of births/deaths.  $dt$  defines the time interval over which this process is defined.

718 The total size of the cells within the  $i$ th clone at time  $t$  as  $X_i(t) = \sum_{l \leq t} \lambda_i(l) - \psi_i(l)$ . A single  
719 parameter  $s_i$  determines the likelihood of a given HSC falling into state 1. or 2., and thus we refer to this  
720 parameter as the clone “fitness.”

721 At any given time  $t$ , the VAF of the  $i$ th clone is defined as  $VAF_i(t) = \frac{X_i(t)}{\sum_j X_j(t)}$ . We define the number of  
722 passenger mutations at time  $t$  in the  $i$ th clone as  $A_i(t) \sim \text{Poisson}(X_i(t) * \mu_p * dt)$ , where  $\mu_p$  is a per-  
723 cell passenger mutation rate. We define  $AC_i(t)$  as the count of “censored” passenger mutations at time  
724  $t$  for the  $i$ th clone, where the censoring occurs due to the limited sensitivity of ~38x sequencing  
725 coverage. This censoring is implemented by the following probability  $P(\text{Binomial}(38, VAF_i) > 2)$ .

726 Association between GEM and gene expression

727 Separate association analyses were performed for each of the five tissue types available within TOPMed:  
728 whole blood, PBMCs, T cells, monocytes, and nasal epithelial tissue. We performed linear regression  
729 between the inverse normalized GEM estimate of the true mutation burden and inverse normalized  
730 gene expression in the tissue, where chronological age, sex, germline genotype PCs 1-15 and expression  
731 PCs 1-20 were included as covariates. In the whole blood analysis, we also included a cohort indicator as  
732 a covariate. Summary statistics from each analysis were then included a Bayesian multivariate analysis  
733 implemented in mashr<sup>62</sup>. As a measure of “significance”, we used the mashr estimate of the local false-  
734 sign rate (LFSR) < 0.05. Enrichment analyses were performed with the pathfinder<sup>80</sup> package including all  
735 tested genes as the background set and Reactome<sup>81</sup> as the reference database for gene sets.

736

#### 737 Incident ischemic stroke analysis

738 Ischemic stroke at most recent visit was chosen for the survival analysis event, and the time to event  
739 was defined as the difference in years between baseline and the most recent visit. The WHI, CHS, and  
740 ARIC cohorts were included. There was a total of 9,885 individuals included in this analysis from the WHI  
741 cohort. In WHI, 9,520 samples were included, there were 1,134 events. In CHS, 2,822 samples were  
742 included, and 199 had events. For ARIC, 3,475 samples were included with 231 events. Covariates  
743 included Ischemic case status at baseline, BMI measured at baseline, "ever smoker" status at baseline, a  
744 spline of age at blood draw, and genetic ancestry PCs 1-4.

#### 745 Incident coronary artery disease analysis

746 Incident coronary artery disease at most recent visit was chosen for the survival analysis event, and the  
747 time to event was defined as the difference in years between baseline and the most recent visit. A  
748 composite coronary artery disease phenotype was defined as an event if at least one of the following  
749 occurred during the follow-up period: myocardial infarction, coronary artery bypass graft, angina,  
750 angioplasty, or death due to coronary heart disease. Individuals with prevalent disease based on this  
751 composite phenotype were excluded. The WHI, CHS, COPDGene, and FHS cohorts were included. In  
752 WHI, 9,039 samples were included, there were 1,787 events. In CHS, 2,456 samples were included, and  
753 933 had events. In FHS, 3,786 samples were included and 525 had events. For COPD, 4,987 samples  
754 were included with 133 events. Covariates included Ischemic case status at baseline, BMI measured at  
755 baseline, "ever smoker" status at baseline, a spline of age at blood draw, and genetic ancestry PCs 1-4.

#### 756 Lentiviral transduction of healthy CD34+ cells

757 Lentiviral vectors expressing NRIP1 (V2LHS\_172503, V2LHS\_172504 and V2LHS\_172507) or SMC4-  
758 targeting shRNA (V2LHS\_21882, V3LHS318029, V3LHS\_318030) (Horizon) or non-silencing pGIPZ-puro  
759 lentiviral vector was transfected together with pCMV-dR8.9 and vesicular stomatitis virus G-expressing  
760 plasmids into HEK 293-FT cells using Lipofectamine 2000 (Thermo Fisher Scientific) for lentiviral  
761 supernatant production as previously described<sup>82</sup>. Primary CD34+ cells were obtained as excess material  
762 from harvests of normal donors for allogeneic bone marrow transplantation. Specimens were collected  
763 by the Johns Hopkins Kimmel Cancer Center Specimen Accessioning Core. Appropriate informed consent  
764 was obtained from all donors before specimen collection in accordance with the Declaration of Helsinki  
765 and under a research protocol approved by the Johns Hopkins Institutional Review Board. CD34+ cell  
766 subsets were isolated using the CD34 MicroBead kit (Miltenyi Biotec) as previously described<sup>83</sup>. CD34+  
767 cells were incubated with the viral supernatant and polybrene (8µg/ml; MilliporeSigma) for transduction  
768 in wells pre-coated with retronectin (20ng/ml; MilliporeSigma). After at least 48 hours, cells were treated  
769 with puromycin (0.5µg/ml; MilliporeSigma) for 4 days to select resistant cells.

#### 770 Apoptosis and differentiation assays

771 Apoptosis was assessed by 7-AAD staining evaluated by flow cytometry (Thermo Fisher Scientific #00-  
772 6993-50). Percentages of stem (CD34+CD38-) and progenitor cells (CD34+CD38+) were assessed by CD34  
773 (Thermo Fisher Scientific #11-0349-42) and CD38 (BioLegend #356641) staining evaluated by flow  
774 cytometry.

775



776 Clonogenicity assays  
777 CD34+ cells following puromycin treatment were collected, counted, and plated at a density of 2000  
778 cells/ml in methylcellulose-based media as previously described<sup>82</sup>. After 10 to 14 days of incubation at  
779 37°C in 5% CO<sub>2</sub>, the recovery of colony-forming units (burst forming unit-erythroid (BFU-E) and colony  
780 forming unit-granulocyte/monocyte (CFU-GM)) were determined by colony counting under bright-field  
781 microscopy. A cell aggregate composed of >50 cells was defined as a colony.

## 782 Code and data availability

783 Code for the Genomic and Epigenomic Mutation rate pipeline: <https://github.com/weinstockj/GEM> .  
784 Individual whole-genome sequence data for TOPMed whole genomes, individual-level harmonized  
785 phenotypes and the CHIP variant call sets used in this analysis are available through restricted access via  
786 the dbGaP TOPMed Exchange Area available to TOPMed investigators.

## 787 Acknowledgements

788 Whole genome sequencing (WGS) for the Trans-Omics in Precision Medicine (TOPMed) program was  
789 supported by the National Heart, Lung and Blood Institute (NHLBI). See Supplementary Information 1  
790 for study omics support information. Centralized read mapping and genotype calling, along with variant  
791 quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-  
792 117626-02S1; contract HHSN268201800002I). Phenotype harmonization, data management, sample-  
793 identity quality control and general study coordination were provided by the TOPMed Data Coordinating  
794 Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). We thank the studies and  
795 participants who provided biological samples and data for TOPMed. The full study-specific  
796 acknowledgments are included in Supplementary Cohort Acknowledgements. The views expressed in  
797 this manuscript are those of the authors and do not necessarily represent the views of the National  
798 Heart, Lung, and Blood Institute; the National Institutes of Health; or the US Department of Health and  
799 Human Services. The authors wish to acknowledge the contributions of the consortium working on the  
800 development of the NHLBI BioData Catalyst ecosystem.

## 801 Competing Interests Declaration

802 L.M.R. is a consultant for the TOPMed Administrative Coordinating Center (through Westat). B.M.P.  
803 serves on the Steering Committee of the Yale Open Data Access Project funded by Johnson & Johnson.  
804 J.Y. reports grant support from Bayer. M.C. reports grant support from Bayer and GSK, Consulting and  
805 speaking fees from Illumina and AstraZeneca. A.G.B. , P.N, and S.J. are cofounders, equity holders, and  
806 on the scientific advisory board of TenSixteen Bio. G.R.A. is an employee of Regeneron  
807 Pharmaceuticals and receives salary, stock and stock options as compensation.

808

## 809 Works Cited

- 810 1. Jaiswal, S. *et al.* Age-Related Clonal Hematopoiesis Associated with Adverse Outcomes A BS TR AC T.  
811 *NEJM.org. N Engl J Med* **26**, 2488–98 (2014).
- 812 2. Genovese, G. *et al.* Clonal Hematopoiesis and Blood-Cancer Risk Inferred from Blood DNA Sequence.  
813 *New England Journal of Medicine* **371**, 2477–2487 (2014).
- 814 3. Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion and  
815 malignancies. *Nature Medicine* **20**, 1472–1478 (2014).
- 816 4. Jaiswal, S. *et al.* Clonal Hematopoiesis and Risk of Atherosclerotic Cardiovascular Disease. *New*  
817 *England Journal of Medicine* (2017) doi:10.1056/NEJMoa1701719.
- 818 5. Desai, P. *et al.* Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nature*  
819 *Medicine* **24**, 1015–1023 (2018).
- 820 6. Bick Alexander G. *et al.* Genetic Interleukin 6 Signaling Deficiency Attenuates Cardiovascular Risk in  
821 Clonal Hematopoiesis. *Circulation* **141**, 124–131 (2020).
- 822 7. Steensma, D. P. *et al.* Clonal hematopoiesis of indeterminate potential and its distinction from  
823 myelodysplastic syndromes. *Blood* **126**, 9–16 (2015).
- 824 8. Loh, P.-R. *et al.* Insights about clonal hematopoiesis from 8,342 mosaic chromosomal alterations.  
825 *Nature* **559**, 350–355 (2018).
- 826 9. Loh, P.-R., Genovese, G. & McCarroll, S. A. Monogenic and polygenic inheritance become  
827 instruments for clonal selection. *Nature* **584**, 136–141 (2020).
- 828 10. Terao, C. *et al.* Chromosomal alterations among age-related haematopoietic clones in Japan. *Nature*  
829 **584**, 130–135 (2020).
- 830 11. Zink, F. *et al.* Clonal hematopoiesis, with and without candidate driver mutations, is common in the  
831 elderly. *Blood* **130**, 742–752 (2017).

- 832 12. Mitchell, E. *et al.* Clonal dynamics of haematopoiesis across the human lifespan. *Nature* **606**, 1–8  
833 (2022).
- 834 13. Weinstock, J. S. *et al.* Aberrant activation of TCL1A promotes stem cell expansion in clonal  
835 haematopoiesis. *Nature* **616**, 755–763 (2023).
- 836 14. Zekavat, S. M. *et al.* Hematopoietic mosaic chromosomal alterations increase the risk for diverse  
837 types of infection. *Nat Med* **27**, 1012–1024 (2021).
- 838 15. Sano, S. *et al.* Hematopoietic loss of Y chromosome leads to cardiac fibrosis and heart failure  
839 mortality. *Science* **377**, 292–297 (2022).
- 840 16. Schuster-Böckler, B. & Lehner, B. Chromatin organization is a major influence on regional mutation  
841 rates in human cancer cells. *Nature* **488**, 504–507 (2012).
- 842 17. Bick, A. G. *et al.* Inherited causes of clonal haematopoiesis in 97,691 whole genomes. *Nature* **586**,  
843 763–768 (2020).
- 844 18. Kar, S. P. *et al.* *Genome-wide analyses of 200,453 individuals yields new insights into the causes and*  
845 *consequences of clonal hematopoiesis.* 2022.01.06.22268846  
846 <https://www.medrxiv.org/content/10.1101/2022.01.06.22268846v1> (2022)  
847 doi:10.1101/2022.01.06.22268846.
- 848 19. Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421–  
849 427 (2014).
- 850 20. Jakubek, Y. A. *et al.* Mosaic chromosomal alterations in blood across ancestries using whole-genome  
851 sequencing. *Nat Genet* 1–8 (2023) doi:10.1038/s41588-023-01553-1.
- 852 21. Thompson, D. J. *et al.* Genetic predisposition to mosaic Y chromosome loss in blood. *Nature* **575**,  
853 652–657 (2019).
- 854 22. Zhou, W. *et al.* Mosaic loss of chromosome Y is associated with common variation near TCL1A.  
855 *Nature Genetics* **48**, 563–568 (2016).

- 856 23. Terao, C. *et al.* GWAS of mosaic loss of chromosome Y highlights genetic effects on blood cell  
857 differentiation. *Nat Commun* **10**, 4719 (2019).
- 858 24. Stacey, S. N. *et al.* Genetics and epidemiology of mutational barcode-defined clonal hematopoiesis.  
859 *Nat Genet* 1–11 (2023) doi:10.1038/s41588-023-01555-z.
- 860 25. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*  
861 **590**, 290–299 (2021).
- 862 26. Kessler, M. D. *et al.* Common and rare variant associations with clonal haematopoiesis phenotypes.  
863 *Nature* 1–9 (2022) doi:10.1038/s41586-022-05448-9.
- 864 27. Bao, E. L. *et al.* Inherited myeloproliferative neoplasm risk affects haematopoietic stem cells. *Nature*  
865 **586**, 769–775 (2020).
- 866 28. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and characterization.  
867 *Nature Methods* **9**, 215–216 (2012).
- 868 29. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes.  
869 *Nature* **518**, 317–330 (2015).
- 870 30. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 871 31. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded 1000 Genomes  
872 Project cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).
- 873 32. Ferraro, N. M. *et al.* Transcriptomic signatures across human tissues identify functional rare genetic  
874 variation. *Science (New York, N.Y.)* **369**, (2020).
- 875 33. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting the  
876 deleteriousness of variants throughout the human genome. *Nucleic Acids Research* **47**, D886–D894  
877 (2019).
- 878 34. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-  
879 scale genetic association studies. *Nature Genetics* **50**, 1335–1341 (2018).

- 880 35. Kessler, M. D. *et al.* Exome sequencing of 628,388 individuals identifies common and rare variant  
881 associations with clonal hematopoiesis phenotypes. 2021.12.29.21268342  
882 <https://www.medrxiv.org/content/10.1101/2021.12.29.21268342v1> (2022)  
883 doi:10.1101/2021.12.29.21268342.
- 884 36. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in  
885 regression, with application to genetic fine mapping. *Journal of the Royal Statistical Society: Series B*  
886 (*Statistical Methodology*) **82**, 1273–1300 (2020).
- 887 37. Carvalho-Silva, D. *et al.* Open Targets Platform: new developments and updates two years on.  
888 *Nucleic Acids Res* **47**, D1056–D1065 (2019).
- 889 38. Fulco, C. P. *et al.* Activity-by-contact model of enhancer–promoter regulation from thousands of  
890 CRISPR perturbations. *Nature Genetics* **51**, 1664–1669 (2019).
- 891 39. Weinstock, J. S. *et al.* Clonal hematopoiesis is driven by aberrant activation of TCL1A.  
892 2021.12.10.471810 Preprint at <https://doi.org/10.1101/2021.12.10.471810> (2021).
- 893 40. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using RegulomeDB.  
894 *Genome Res* **22**, 1790–1797 (2012).
- 895 41. Vaser, R., Adusumalli, S., Leng, S. N., Sikic, M. & Ng, P. C. SIFT missense predictions for genomes. *Nat*  
896 *Protoc* **11**, 1–9 (2016).
- 897 42. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456  
898 humans. *Nature* **581**, 434–443 (2020).
- 899 43. Petryszak, R. *et al.* Expression Atlas update—an integrated database of gene and protein expression  
900 in humans, animals and plants. *Nucleic Acids Res* **44**, D746–D752 (2016).
- 901 44. Toren, A. *et al.* CD133-Positive Hematopoietic Stem Cell “Stemness” Genes Contain Many Genes  
902 Mutated or Abnormally Expressed in Leukemia. *Stem Cells* **23**, 1142–1153 (2005).

- 903 45. Chen, M.-H. *et al.* Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667 Individuals from  
904 5 Global Populations. *Cell* **182**, 1198-1213.e14 (2020).
- 905 46. Yuan, R. *et al.* Genetic coregulation of age of female sexual maturation and lifespan through  
906 circulating IGF1 among inbred mouse strains. *Proceedings of the National Academy of Sciences* **109**,  
907 8224–8229 (2012).
- 908 47. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature*  
909 **622**, 329–338 (2023).
- 910 48. Blake, J. A. *et al.* Mouse Genome Database (MGD): Knowledgebase for mouse-human comparative  
911 biology. *Nucleic Acids Res* **49**, D981–D987 (2021).
- 912 49. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nat Genet* **25**, 25–29 (2000).
- 913 50. Vösa, U. *et al.* Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and  
914 polygenic scores that regulate blood gene expression. *Nat Genet* **53**, 1300–1310 (2021).
- 915 51. Fairfax, B. P. *et al.* Innate Immune Activity Conditions the Effect of Regulatory Variants upon  
916 Monocyte Gene Expression. *Science* **343**, 1246949 (2014).
- 917 52. Dhindsa, R. S. *et al.* Rare variant associations with plasma protein levels in the UK Biobank. *Nature*  
918 **622**, 339–347 (2023).
- 919 53. Ye, H. *et al.* Leukemic Stem Cells Evade Chemotherapy by Metabolic Adaptation to an Adipose  
920 Tissue Niche. *Cell Stem Cell* **19**, 23–37 (2016).
- 921 54. Li, X. *et al.* Dynamic incorporation of multiple in silico functional annotations empowers rare variant  
922 association analysis of large whole-genome sequencing studies at scale. *Nat Genet* **52**, 969–983  
923 (2020).
- 924 55. Yeung, Y. T. *et al.* CELF2 suppresses non-small cell lung carcinoma growth by inhibiting the PREX2-  
925 PTEN interaction. *Carcinogenesis* **41**, 377–389 (2020).

- 926 56. Li, Z. *et al.* Dynamic Scan Procedure for Detecting Rare-Variant Association Regions in Whole-  
927 Genome Sequencing Studies. *The American Journal of Human Genetics* **104**, 802–814 (2019).
- 928 57. Chen, L. *et al.* Genetic Drivers of Epigenetic and Transcriptional Variation in Human Immune Cells.  
929 *Cell* **167**, 1398-1414.e24 (2016).
- 930 58. Denny, J. C. *et al.* PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-  
931 disease associations. *Bioinformatics* **26**, 1205–1210 (2010).
- 932 59. Bycroft, C. *et al.* Genome-wide genetic data on ~500,000 UK Biobank participants. *bioRxiv* 166298–  
933 166298 (2017) doi:10.1101/166298.
- 934 60. FinnGen. FinnGen. *FinnGen Documentation of R3 release*  
935 <https://finngen.gitbook.io/documentation/>.
- 936 61. Oliva, M. *et al.* The impact of sex on gene expression across human tissues. *Science* **369**, eaba3066  
937 (2020).
- 938 62. Urbut, S. M., Wang, G., Carbonetto, P. & Stephens, M. Flexible statistical methods for estimating  
939 and testing effects in genomic studies with multiple conditions. *Nat Genet* **51**, 187–195 (2019).
- 940 63. Stephens, M. False discovery rates: a new deal. *Biostatistics* **18**, 275–294 (2017).
- 941 64. Valette, K. *et al.* Prioritization of candidate causal genes for asthma in susceptibility loci derived  
942 from UK Biobank. *Commun Biol* **4**, 700 (2021).
- 943 65. Kanehisa, M., Furumichi, M., Sato, Y., Kawashima, M. & Ishiguro-Watanabe, M. KEGG for taxonomy-  
944 based analysis of pathways and genomes. *Nucleic Acids Res* **51**, D587–D592 (2023).
- 945 66. Anguille, S. *et al.* Interferon- $\alpha$  in acute myeloid leukemia: an old drug revisited. *Leukemia* **25**, 739–  
946 748 (2011).
- 947 67. de Almeida, P. E. *et al.* Anti-VEGF Treatment Enhances CD8+ T-cell Antitumor Activity by Amplifying  
948 Hypoxia. *Cancer Immunol Res* **8**, 806–818 (2020).

- 949 68. Palazon, A. *et al.* An HIF-1 $\alpha$ /VEGF-A Axis in Cytotoxic T Cells Regulates Tumor Progression. *Cancer*  
950 *Cell* **32**, 669-683.e5 (2017).
- 951 69. Heyde, A. *et al.* Increased stem cell proliferation in atherosclerosis accelerates clonal hematopoiesis.  
952 *Cell* **184**, 1348-1361.e22 (2021).
- 953 70. Yu, B. *et al.* Association of Clonal Hematopoiesis With Incident Heart Failure. *Journal of the*  
954 *American College of Cardiology* **78**, 42–52 (2021).
- 955 71. Yu, Z. *et al.* Human Plasma Proteomic Profile of Clonal Hematopoiesis. 2023.07.25.550557 Preprint  
956 at <https://doi.org/10.1101/2023.07.25.550557> (2023).
- 957 72. Regier, A. A. *et al.* Functional equivalence of genome sequencing analysis pipelines enables  
958 harmonized variant calling across human genetics projects. *Nature Communications* **9**, 1–8 (2018).
- 959 73. Jun, G., Wing, M. K., Abecasis, G. R. & Kang, H. M. An efficient and scalable analysis framework for  
960 variant extraction and refinement from population-scale DNA sequence data. *Genome Res* **25**, 918–  
961 925 (2015).
- 962 74. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous  
963 cancer samples. *Nature Biotechnology* **31**, 213–219 (2013).
- 964 75. Voss, K., Gentry, J. & Van der Auwera, G. Full-stack genomics pipelining with GATK4 + WDL +  
965 Cromwell. in (F1000 Research, 2017). doi:10.7490/f1000research.1114631.1.
- 966 76. Beauchamp, E. M. *et al.* ZBTB33 Is Mutated in Clonal Hematopoiesis and Myelodysplastic  
967 Syndromes and Impacts RNA Splicing. *Blood Cancer Discov* (2021) doi:10.1158/2643-3230.BCD-20-  
968 0224.
- 969 77. Weinstock, J. S. *et al.* The genetic determinants of recurrent somatic mutations in 43,693 blood  
970 genomes. *Science Advances* **9**, eabm4945 (2023).
- 971 78. Jun, G. *et al.* Structural variation across 138,134 samples in the TOPMed consortium. *bioRxiv*  
972 2023.01.25.525428 (2023) doi:10.1101/2023.01.25.525428.



- 973 79. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biology* **17**, 122 (2016).
- 974 80. Ulgen, E., Ozisik, O. & Sezerman, O. U. pathfindR: An R Package for Comprehensive Identification of  
975 Enriched Pathways in Omics Data Through Active Subnetworks. *Frontiers in Genetics* **10**, (2019).
- 976 81. Gillespie, M. *et al.* The reactome pathway knowledgebase 2022. *Nucleic Acids Research* **50**, D687–  
977 D692 (2022).
- 978 82. Karantanos, T. *et al.* The role of the atypical chemokine receptor CCRL2 in myelodysplastic  
979 syndrome and secondary acute myeloid leukemia. *Sci Adv* **8**, eabl8952 (2022).
- 980 83. Karantanos, T. *et al.* CCRL2 affects the sensitivity of myelodysplastic syndrome and secondary acute  
981 myeloid leukemia cells to azacitidine. *Haematologica* **108**, 1886–1899 (2023).
- 982