

UCLA

UCLA Previously Published Works

Title

The genetic determinants of recurrent somatic mutations in 43,693 blood genomes

Permalink

<https://escholarship.org/uc/item/1496g128>

Journal

Science Advances, 9(17)

ISSN

2375-2548

Authors

Weinstock, Joshua S
Laurie, Cecelia A
Broome, Jai G
[et al.](#)

Publication Date

2023-04-28

DOI

10.1126/sciadv.abm4945

Peer reviewed

GENETICS

The genetic determinants of recurrent somatic mutations in 43,693 blood genomes

Joshua S. Weinstock^{1*}, Cecelia A. Laurie², Jai G. Broome^{2,3}, Kent D. Taylor⁴, Xiuqing Guo⁴, Alan R. Shuldiner⁵, Jeffrey R. O'Connell⁵, Joshua P. Lewis⁵, Eric Boerwinkle⁶, Kathleen C. Barnes⁷, Nathalie Chami^{8,9}, Eimear E. Kenny¹⁰, Ruth J. F. Loos^{8,9}, Myriam Fornage¹¹, Susan Redline^{12,13}, Brian E. Cade^{12,13,14}, Frank D. Gilliland¹⁵, Zhanghua Chen¹⁵, W. James Gauderman¹⁵, Rajesh Kumar^{16,17}, Leslie Grammer¹⁷, Robert P. Schleimer¹⁷, Bruce M. Psaty^{18,19,20}, Joshua C. Bis¹⁸, Jennifer A. Brody¹⁸, Edwin K. Silverman²¹, Jeong H. Yun²¹, Dandi Qiao^{21†}, Scott T. Weiss^{22,12}, Jessica Lasky-Su^{22,12}, Dawn L. DeMeo^{22,12}, Nicholette D. Palmer²³, Barry I. Freedman²⁴, Donald W. Bowden²³, Michael H. Cho²⁵, Ramachandran S. Vasani²⁶, Andrew D. Johnson^{26,27}, Lisa R. Yanek²⁸, Lewis C. Becker²⁸, Sharon Kardia²⁹, Jiang He³⁰, Robert Kaplan³¹, Susan R. Heckbert^{32,33}, Nicholas L. Smith^{32,33,34}, Kerri L. Wiggins³⁵, Donna K. Arnett³⁶, Marguerite R. Irvin³⁷, Hemant Tiwari³⁷, Adolfo Correa³⁸, Laura M. Raffield³⁹, Yan Gao⁴⁰, Mariza de Andrade⁴¹, Jerome I. Rotter⁴, Stephen S. Rich⁴², Ani W. Manichaikul⁴², Barbara A. Konkle⁴³, Jill M. Johnsen^{44,43}, Marsha M. Wheeler⁴⁵, Brian S. Custer⁴⁶, Ravindranath Duggirala^{47,48}, Joanne E. Curran^{47,48}, John Blangero^{47,48}, Hongsheng Gui^{49,50}, Shujie Xiao^{49,50}, L. Keoki Williams^{49,50}, Deborah A. Meyers⁵¹, Xingnan Li⁵², Victor Ortega⁵³, Stephen McGarvey⁵⁴, C. Charles Gu⁵⁵, Yii-Der Ida Chen⁴, Wen-Jane Lee⁵⁶, M. Benjamin Shoemaker⁵⁷, Dawood Darbar⁵⁸, Dan Roden⁵⁹, Christine Albert⁶⁰, Charles Kooperberg⁶¹, Pinkal Desai^{62,63}, Thomas W. Blackwell¹, Goncalo R. Abecasis^{1,64}, Albert V. Smith¹, Hyun M. Kang¹, Rasika Mathias²⁸, Pradeep Natarajan^{14,65,66}, Siddhartha Jaiswal⁶⁷, Alexander P. Reiner^{61,68}, Alexander G. Bick^{69*}, NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium

Nononcogenic somatic mutations are thought to be uncommon and inconsequential. To test this, we analyzed 43,693 National Heart, Lung and Blood Institute Trans-Omics for Precision Medicine blood whole genomes from 37 cohorts and identified 7131 non-missense somatic mutations that are recurrently mutated in at least 50 individuals. These recurrent non-missense somatic mutations (RNMSMs) are not clearly explained by other clonal phenomena such as clonal hematopoiesis. RNMSM prevalence increased with age, with an average 50-year-old having 27 RNMSMs. Inherited germline variation associated with RNMSM acquisition. These variants were found in genes involved in adaptive immune function, proinflammatory cytokine production, and lymphoid lineage commitment. In addition, the presence of eight specific RNMSMs associated with blood cell traits at effect sizes comparable to Mendelian genetic mutations. Overall, we found that somatic mutations in blood are an unexpectedly common phenomenon with ancestry-specific determinants and human health consequences.

INTRODUCTION

As humans age, their dividing cells acquire mutations. Except for the rare cases when a mutation occurs in a driver gene—possibly leading to cancer (*1*)—the vast majority of these somatic mutations do not alter protein coding sequence. Therefore, they are believed to have little effect on cellular function. However, this assumption has not been rigorously tested and may not be true, as noncoding regions can actively regulate gene expression through enhancers, splicing, genomic structure, and chromatin configuration.

If nononcogenic somatic mutations affect physiology, then the blood cells would be one of the cell types most affected. Erythrocytes and immune cells are constantly regenerated by hematopoietic stem cells (HSCs), which have accumulated a lifetime of mutations, and these blood cells are responsible for distributing oxygen, mediating immune responses, and regulating clotting—all of which could be affected by somatic mutations. Studies have found seemingly non-cancerous clonal expansion of HSCs in healthy individuals,

suggesting that somatic mutation is common (2–8). They also found that individuals with clonal hematopoiesis have increased risk of not only blood cancer but also an array of nonmalignant aging diseases, including heart disease, suggesting that somatic mutations can affect physiology. However, these studies only assessed coding variation. Recent efforts to identify noncoding somatic mutations have been performed only in a limited number of people, such as the recent Pan-Cancer Analysis of Whole Genomes (PCAWG) analysis, which included only 2658 genomes (9).

Here, using 43,693 deeply sequenced (38×) whole genomes from the National Heart, Lung and Blood Institute Trans-Omics for Precision Medicine (TOPMed) initiative (*10*), we developed a catalog of somatic mutations derived from peripheral blood samples of individuals from the general population and identify a class of recurring mutations across individuals, which we term recurrent non-missense somatic mutations (RNMSMs). We then conducted the first

Copyright © 2023 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works. Distributed under a Creative Commons Attribution NonCommercial License 4.0 (CC BY-NC).

large-scale investigation of their germline determinants and health consequences.

RESULTS

Identifying recurrent somatic mutations

In total, we used the whole-genome data of 47,243 individuals from 37 cohorts included in TOPMed, which were sequenced at seven distinct genome sequencing centers (tables S1 to S3), and identified somatic point mutations in their peripheral blood by running GATK Mutect2 (11). We applied stringent variant filtering procedures to deplete our call set of sequencing artifacts and germline variants (Fig. 1A, Methods, and text S1). We also excluded all calls with a variant allele fraction (VAF) of >35%, defined as the percentage of reads that are alternative at a given site. Of the 14,841,388 somatic variants identified, most (89.9%) were only observed in one individual (Fig. 1B). These singleton variants are likely passenger mutations on a clonally expanded HSC because, at the read depth sequenced, it is unlikely that single variants would be detected unless they were present in a large enough proportion of peripheral blood cells (5, 12).

Intprisingly, we identified a small subset of somatic genetic variants that are recurrently mutated and were neither missense nor

nonsense mutations, which we call RNMSMs (synonymous mutations are retained). We then further refined our RNMSM calls with additional germline and sequencing artifact filters. We first excluded high-quality germline variants that were identified in the TOPMed germline single-nucleotide polymorphism (SNP) calls (13). Next, using the TOPMed germline structural variant calls (14), we also excluded all mutations that overlapped with a common [minor allele fraction (MAF) > 10%] germline duplication, as such events may result in artifactual somatic variant calls. As sequencing artifacts may have similar signatures to somatic mutations, we also excluded parts of the genome that were not uniquely mappable, and we excluded mutations with mapping artifact signatures (text S1). We then excluded samples where age at blood draw was not available, resulting in an analysis set of 43,693 samples. We focused our analysis on the RNMSMs that occurred at least 50 times and in at least five distinct cohorts to deplete for cohort-specific sequencing artifacts. This yielded a set of 7131 RNMSMs (table S4) or ~0.05% of all variants. The average number of RNMSMs per sample was 28.6, with a median of 25, a minimum of 1, a maximum of 163, and an SD of 16.1.

The RNMSM burden was positively associated with age at the time of blood draw [in units of decades, $\beta = 1.4$, 95% confidence

¹Center for Statistical Genetics, Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, MI 48109, USA. ²Department of Biostatistics, University of Washington, Seattle, WA 98195, USA. ³Division of Medical Genetics, Department of Medicine, University of Washington, Seattle, WA 98195, USA. ⁴The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA 90502, USA. ⁵Department of Medicine, University of Maryland, Baltimore, Baltimore, MD 21201, USA. ⁶Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA. ⁷Division of Biomedical Informatics and Personalized Medicine, Department of Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO 80045, USA. ⁸The Charles Bronfman Institute of Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. ⁹The Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. ¹⁰Institute for Genomic Health, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. ¹¹Brown Foundation Institute of Molecular Medicine, McGovern Medical School, University of Texas Health Science Center at Houston, Houston, TX 77030, USA. ¹²Department of Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA. ¹³Harvard Medical School, Boston, MA 02115, USA. ¹⁴Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA. ¹⁵Department of Preventive Medicine, University of Southern California, Los Angeles, CA 90089, USA. ¹⁶Ann and Robert H. Lurie Children's Hospital of Chicago, Chicago, IL 60611, USA. ¹⁷Northwestern University Feinberg School of Medicine, Chicago, IL 60611, USA. ¹⁸Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA 98195, USA. ¹⁹Department of Epidemiology, University of Washington, Seattle, WA 98195, USA. ²⁰Department of Medicine, University of Washington, Seattle, WA 98195, USA. ²¹Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA. ²²Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA. ²³Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC 27101, USA. ²⁴Department of Internal Medicine, Section on Nephrology, Wake Forest School of Medicine, Winston-Salem, NC 27101, USA. ²⁵Channing Division of Network Medicine and Division of Pulmonary and Critical Care Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA. ²⁶National Heart, Lung, and Blood Institute's, Boston University's Framingham Heart Study, Framingham, MA 01701, USA. ²⁷National Heart, Lung and Blood Institute, Population Sciences Branch, Framingham, MA 01701, USA. ²⁸Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA. ²⁹Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI 48109, USA. ³⁰Department of Epidemiology, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA 70112, USA. ³¹Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY 10461, USA. ³²Department of Epidemiology, University of Washington, Seattle, WA 98195, USA. ³³Kaiser Permanente Washington Health Research Institute, Kaiser Permanente Washington, Seattle, WA 98101, USA. ³⁴Seattle Epidemiologic Research and Information Center, Department of Veterans Affairs Office of Research and Development, Seattle, WA 98108, USA. ³⁵Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA 98101, USA. ³⁶Dean's Office, College of Public Health, University of Kentucky, Lexington, KY 40506, USA. ³⁷University of Alabama at Birmingham, Birmingham, AL 35294, USA. ³⁸Department of Medicine, Jackson Heart Study, University of Mississippi Medical Center, Jackson, MS 39216, USA. ³⁹Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA. ⁴⁰Department of Medicine, University of Mississippi Medical Center, Jackson, MS 39216, USA. ⁴¹Department of Health Sciences Research, Mayo Clinic, Rochester, MN 55905, USA. ⁴²Department of Public Health Sciences, Center for Public Health Genomics, University of Virginia, Charlottesville, VA 22903, USA. ⁴³Department of Medicine, University of Washington, Seattle, WA 98195, USA. ⁴⁴Research Institute, Bloodworks Northwest, Seattle, WA 98102, USA. ⁴⁵Genome Science, University of Washington, Seattle, WA 98195, USA. ⁴⁶Vitalant Research Institute, San Francisco, CA 94105, USA. ⁴⁷Department of Human Genetics, University of Texas Rio Grande Valley School of Medicine, Brownsville, TX 78520, USA. ⁴⁸South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley School of Medicine, Brownsville, TX 78520, USA. ⁴⁹Center for Individualized and Genomic Medicine Research (CIGMA), Henry Ford Health System, Detroit, MI 48202, USA. ⁵⁰Department of Medicine, Henry Ford Health System, Detroit, MI 48202, USA. ⁵¹Division of Genetics, Genomics, and Precision Medicine, University of Arizona, Tucson, AZ 85721, USA. ⁵²Department of Medicine, University of Arizona, Tucson, AZ 85721, USA. ⁵³Wake Forest University School of Medicine, Winston-Salem, NC 27101, USA. ⁵⁴Department of Epidemiology and International Health Institute, Brown University School of Public Health, Providence, RI 02903, USA. ⁵⁵Division of Biostatistics, Washington University School of Medicine, Campus Box 8067, 660 S. Euclid Avenue, St. Louis, MO 63110, USA. ⁵⁶Department of Medical Research, Taichung Veterans General Hospital, 1650, Sec. 4, Taiwan Boulevard, Taichung City, Taiwan. ⁵⁷Division of Cardiology, Vanderbilt University Medical Center, Nashville, TN 37232, USA. ⁵⁸Division of Cardiology, University of Illinois at Chicago, Chicago, IL 60607, USA. ⁵⁹Departments of Medicine, Pharmacology, and Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN 37232, USA. ⁶⁰Department of Cardiology, Cedars-Sinai, Los Angeles, CA 90048, USA. ⁶¹Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA. ⁶²Division of Hematology and Oncology, Weill Cornell Medicine, New York, NY 10065, USA. ⁶³Englander Institute of Precision Medicine, Weill Cornell Medicine, New York 10065, NY, USA. ⁶⁴Regeneron Pharmaceuticals, Tarrytown, NY 10591, USA. ⁶⁵Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA 02114, USA. ⁶⁶Department of Medicine, Harvard Medical School, Boston, MA 02115, USA. ⁶⁷Department of Pathology, Stanford University, Stanford, CA 94305, USA. ⁶⁸Department of Epidemiology, University of Washington, Seattle, WA 98195, USA. ⁶⁹Division of Genetic Medicine, Department of Medicine, Vanderbilt University, Nashville, TN 37232, USA.

†Present address: NuvoAir U.S. Inc., Boston, MA 02109, USA.

*Corresponding author. Email: jweinstk@umich.edu (J.S.W.); alexander.bick@vumc.org (A.G.B.)

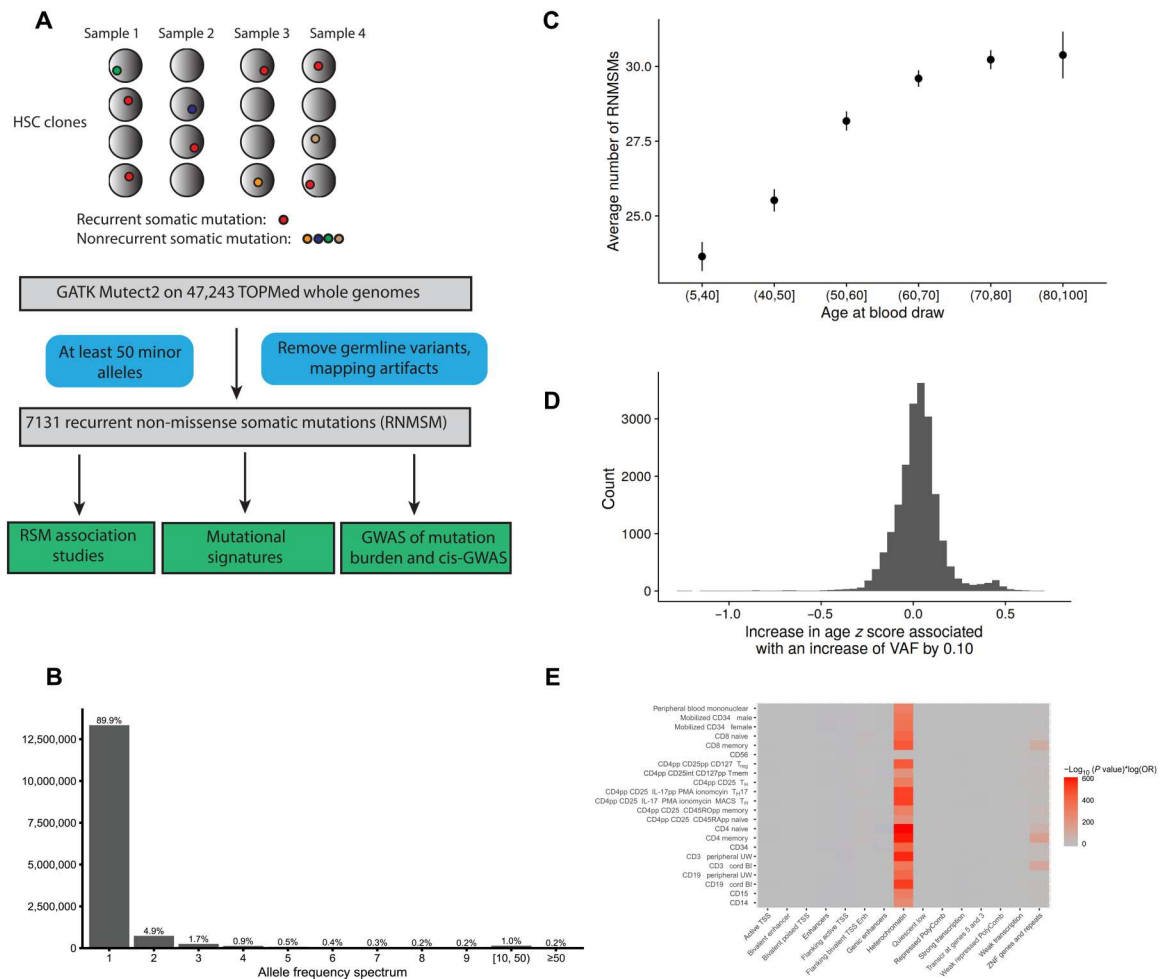


Fig. 1. Identification of recurrent somatic mutations from TOPMed whole genomes. (A) Identification of RNMSMs. (B) The allele frequency spectrum of somatic mutations. These represent the frequency of sets of variants with the given allele count. (C) The average number of RNMSMs stratified across age bins with 95% CIs. (D) The histogram of the slopes of the individual RNMSMs with age. (E) Enrichments of chromHMM annotations across blood epigenomes from Roadmap Epigenomics. Enrichments are defined as the log odds ratio (OR) multiplied by the $-\log_{10}(P \text{ value})$ from the Fisher's exact test. T_{reg} , regulatory T cell; TH, T helper cell; IL-17pp, interleukin-17pp; PMA, phorbol 12-myristate 13-acetate; MACS, magnetic-activated cell sorting; Tmem, T memory.

interval (CI) (1.3, 1.5), $P \text{ value} = 5.3 \times 10^{-135}$; Fig. 1, C and D] likely because older people have had more time to accumulate mutations. The association between RNMSM burden and age was attenuated by the inclusion of a cohort indicator covariate [in units of decades, $\beta = 0.20$, 95% CI (0.08, 0.32), $P \text{ value} = 8.4 \times 10^{-4}$], but it is still unlikely to occur by chance.

We then experimentally validated a subset of the RNMSMs by calling somatic variants from 46 individuals not sequenced as part of TOPMed using 30x whole-genome sequencing (WGS) of paired blood and heart tissue (Methods). We observed that 219 of the RNMSMs detected in the TOPMed study were present in the 46 blood samples (table S5), with at least one alt-allele present in the data. We observed that 266 mutations were present, for an average of 5.8 mutations per sample. The low number of replicated RNMSMs likely reflects the younger age of the validation samples (mean of 50.6 years old in validation and mean of 61.3 in discovery cohorts), the decreased diversity compared to TOPMed, and the lower sequencing depth.

We then further validated the RNMSMs by estimating their Mendelian concordance. We reasoned that RNMSMs should have high rates of Mendelian discordance among trios because they are not germline variants. Using 4103 trios in TOPMed, we calculated the Mendelian genotype discordance among Mendelian-informative meiosis. We then categorized RNMSMs as Mendelian discordant if there were at least two Mendelian discordances present and the genotype discordance was greater than 2%. We observed that 71.4% ($n = 5094$) of the RNMSMs were Mendelian discordant (table S6).

To assess whether the RNMSMs were spurious signatures of rare germline duplications, we used the TOPMed structural variant calls to estimate the co-occurrence of the individual RNMSM genotypes with overlapping germline duplications and deletions. We observed that 9% ($n = 616$) of the RNMSMs co-occurred frequently (>80% of mutations) with overlapping germline duplications and deletions (table S7), suggesting that most RNMSMs are unlikely to tag germline structural variants.

To test whether the RNMSMs could be identified in another population-scale sequencing cohort, we cross-referenced the RNMSMs with the site list of the gnomAD 3.1.2 release (15). Because samples overlap between TOPMed and gnomAD, we used the non-TOPMed allele frequencies reported in the site list. We observed that 6386 of the RNMSMs were reported in the gnomAD site list (table S8). Collectively, these observations highlight that the RNMSMs have been depleted of germline variants and sequencing artifacts and can be detected in other cohorts.

To characterize the function of the RNMSMs, we annotated them using Variant Effect Predictor (VEP) (16). We observed that the RNMSMs were enriched in intergenic regions (twofold increase; fig. S1). We then calculated the enrichment of RNMSMs in cell type-specific chromatin signatures using the Roadmap Epigenomics catalog (17) (see Methods) and found strong enrichment in heterochromatin across blood cell types (Fig. 1E). This is consistent with a previous study that found that 40% of variance in cancer mutation density can be attributed to variation in H3K9me3, a histone mark associated with heterochromatin (18). This association may be mediated by diminished DNA repair activity in heterochromatin (19).

To identify potential causes of RNMSMs, we conducted a mutational signature analysis (see Methods) and found that the mutational signature SBS5 was the only contributing factor. SBS5 is the result of a clock-like mutational process, meaning that the burden of SBS5 mutations correlates with age and may be associated with tobacco smoking and nucleotide excision repair deficiency (20).

We then asked whether the RNMSMs could be explained by clonal hematopoiesis of indeterminate potential (CHIP) (21) or were an independent phenomenon. We stratified the distribution of RNMSMs based on CHIP carrier status (8). We observed a very modest increase in RNMSM number among the CHIP carriers [fig. S2; $\beta = 0.84$, 95% CI (0.43, 1.25), P value = 5.6×10^{-5}], suggesting that RNMSMs are not only explained by clonal driver mutation phenomena.

We observed a small number of RNMSMs in genes that are mutated in clonal hematopoiesis and hematological malignancy, including intronic RNMSMs in *TET2* (rs1341404863, rs1428954790), *JAK2* (rs1223087352), and *ASXL2* (rs1416952509). Among acute myeloid leukemia (AML) genes, we observed a variant downstream of *BCL6* (rs9878379) with 194 mutations and an intronic variant of *FLT3* (rs1242299801). rs9878379 was previously reported in the International Cancer Genome Consortium (ICGC) cancer cohort in six cancer types (22) but has previously been unreported in a non-cancer population cohort. *BCL6* is a proto-oncogene that is a frequently mutated in B cell lymphomas (23).

To estimate the noncancer population prevalence of RNMSMs that occur in cancer genomes, we cross-referenced the RNMSMs with the cancer genome somatic calls released by the PCAWG (24). We observed that 81 of the somatic mutations called in the PCAWG cancer genomes were present among the RNMSMs (table S9). These include (rs1379713562, COSV71733298), a synonymous mutation *SBSN* that was mutated in 191 genomes, and an intronic variant (rs1473587176) in *PTER* with a notable allele frequency (930 mutated alleles, 1% MAF). *SBSN* expression has recently been described as a potential biomarker in myelodysplastic syndromes (25). We observed that 23 of the PCAWG-RNMSMs had at least 100 mutations in our cohort, indicating that a subset

of mutations in cancer genomes is prevalent in a population cohort that was not ascertained for cancer.

RNMSMs and genetic ancestry

We then performed principal components analysis (PCA) on the somatic variant matrix of VAFs to characterize the underlying structure of RNMSMs (Methods). We estimated the first 50 components (26) and observed that the first three axes of variation explain >92% of the variance, indicating strong shared underlying structure among RNMSMs. We observed that the first somatic principal component (sPC1) was strongly associated with overall RNMSM burden ($R_{sq} = 94\%$, P value $< 2 \times 10^{-16}$). As PCs of germline variants associate with genetic ancestry (27–29), we then asked whether sPCs had similar properties. As previously described, the global genetic ancestry of TOPMed samples has been estimated with RFMix (13). Using this resource, we labeled samples according to their largest global ancestry contribution. We then conducted a multinomial logistic regression of global ancestry on the first five sPCs, including a study indicator as a covariate. We observed that sPCs were strongly associated with ancestry (Fig. 2, A and B, and table S10). We then calculated fixation index (F_{st}) values (Methods) among the somatic variants to estimate the degree of population differentiation at individual mutations. We observed that 55 of the somatic mutations had F_{st} estimates > 0.05 (Fig. 2C and table S11). These ancestry associations suggested that RNMSMs may associate with germline variants common in particular ancestries. This is consistent with previously known associations between genetic ancestry and DNA replication timing (30)—DNA recombination hotspots (31). The association between clonal phenomena and genetic ancestry has been previously observed in both CHIP (8) and mosaic chromosomal alterations (32). Therefore, we sought to establish whether there were common germline genetic variants associated with acquiring RNMSMs.

Germline genetic basis of RNMSMs

We performed a multi-ancestry genome-wide association study (GWAS) of the RNMSM burden using the scalable and accurate implementation of generalized mixed model (SAIGE) (33) and identified five genome-wide significant loci, all of which collectively highlight the influence of immune function on RNMSM burden (Fig. 3A). The strongest signal was observed in the *HLA* region at rs9271735, which is a common (MAF 28%) variant and is 2.5 kb upstream of the transcription start site (TSS) of *HLA-DQAI*, whose protein plays an essential role in antigen presentation. The A allele was associated with a 0.09 SD increase in RNMSM burden (P value = 3.8×10^{-98}). Given that there is extensive linkage disequilibrium at the *HLA* locus, this variant likely tags a specific *HLA* haplotype. To ensure that this signal was not a consequence of population stratification, we then performed both European and African ancestry-specific GWAS (Methods; fig. S3). We observed that rs9271735 was genome-wide significant in both ancestry-specific GWAS (African ancestry, P value = 4.9×10^{-18} ; European ancestry, P value = 5.1×10^{-40}), indicating that the association is unlikely to be a consequence of population stratification. Consistent with the possible role of the adaptive immune system in surveillance of HSCs for excessive mutation, a recent report showed that HSCs in humans are antigen-presenting cells (34).

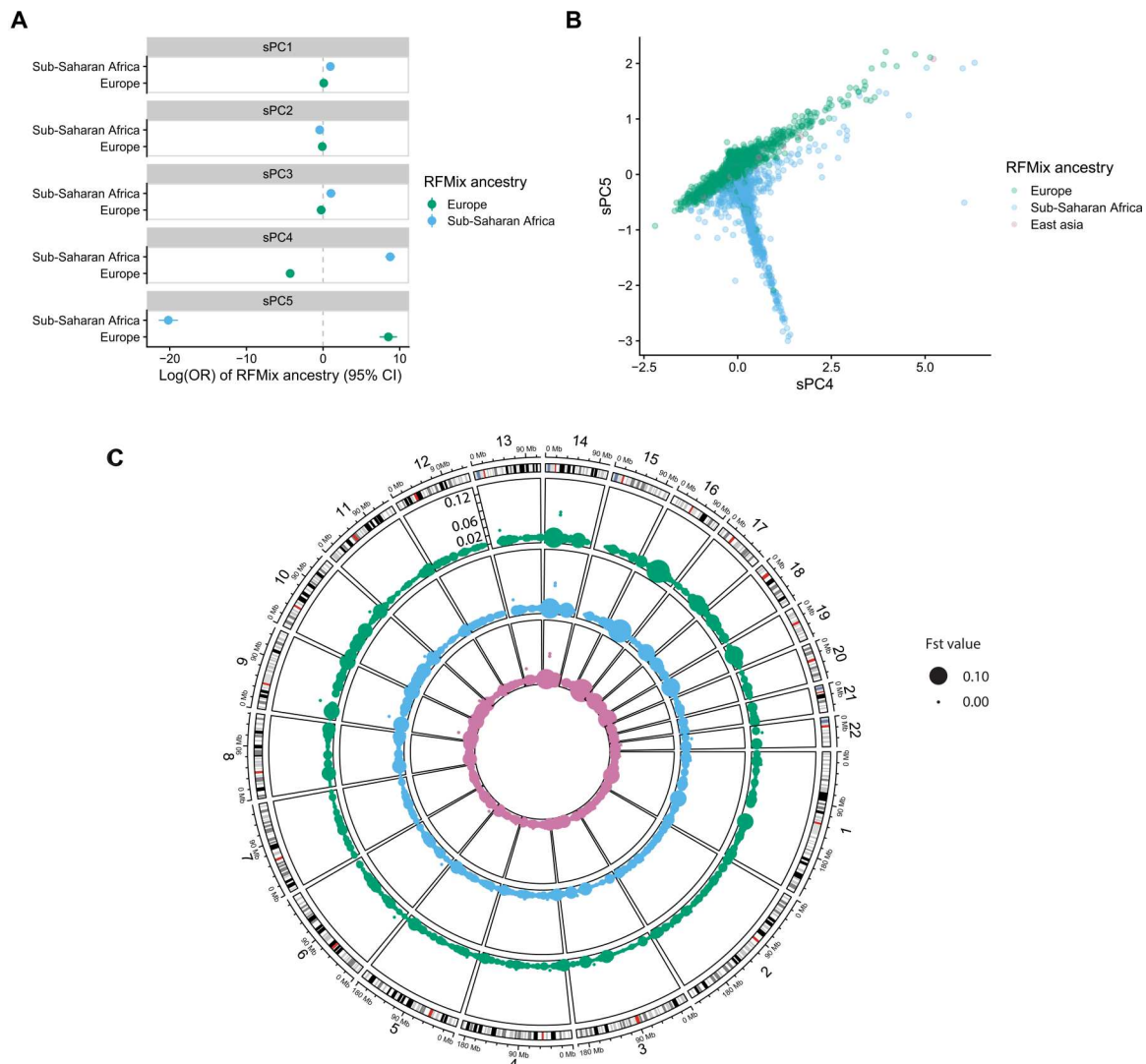


Fig. 2. Somatic genotype PCs are associated with genetic ancestry. (A) Global ancestry labels, as estimated from RFMix, are regressed on the first five sPCs in a multinomial regression with a study indicator included as a covariate. The ORs are estimated with East Asian ancestry as the reference level. (B) A scatterplot of sPC4 and sPC5, with colors indicated by the RFMix global ancestry label. (C) A circular genome plot where the angle indicates genomic position, and the radius within a given track indicates the allele frequency, which ranges from 0 to 0.12. A separate track is plotted for the allele frequencies computed separately in European, Sub-Saharan African, and East Asian genomes, respectively, which are colored by the legend indicated in (B). The size of points indicates the fixation index (Fst) estimate, where larger points have larger Fst values.

The *MIF-GSTT2B* locus was another hit, with the leading variant being rs5760124, a common (MAF 43%) intergenic variant 31.3 and 35.5 kb from the TSSs of *GSTT2B* and *MIF*, respectively. The T allele was associated with a 0.04 SD increase in RNMSM burden (P value = 1.8×10^{-25}). It is an expression quantitative trait locus (eQTL) in whole blood (35) for *MIF*, *DDT*, and *DDTL*, and the T allele increases expression for all three (normalized effect sizes of 0.80, 0.40, and 0.25; P values of 2.6×10^{-60} , 3.5×10^{-55} , and 3.6×10^{-14}). It is also a protein quantitative trait locus (pQTL) for *GSTT2B* in blood plasma (36) and a splicing quantitative trait locus (sQTL) for *DDT* in whole blood (35). Consistent with the role of rs5760124 in altered hematopoiesis, the T allele is also associated with increased neutrophil counts (37). *MIF* is a proinflammatory cytokine that regulates macrophage function (38), with *DDT* and *DDTL* being its functional

homologs. *MIF* has been shown to inhibit p53 in murine models (39), facilitating cell proliferation. Increased expression of *MIF* promotes inflammation, which may result in accelerated aging of HSCs (40), providing fertile substrate for the acquisition of aging-associated clonal phenomena.

We then asked whether genetic variation associates with individual RNMSMs rather than RNMSM burden. Cis-acting genetic variation has been previously reported to associate with somatic mutations including *JAK2* V617F (41–43) and somatic chromosomal mosaicism (44, 45). We focused on the 82 most frequent RNMSMs and performed association analyses between each germline variant in a 2-Mb region surrounding each RNMSM. To identify trans RNMSM regulators, we also tested each of the RNMSM cis regions with all the other RNMSMs. In total, we examined 6724

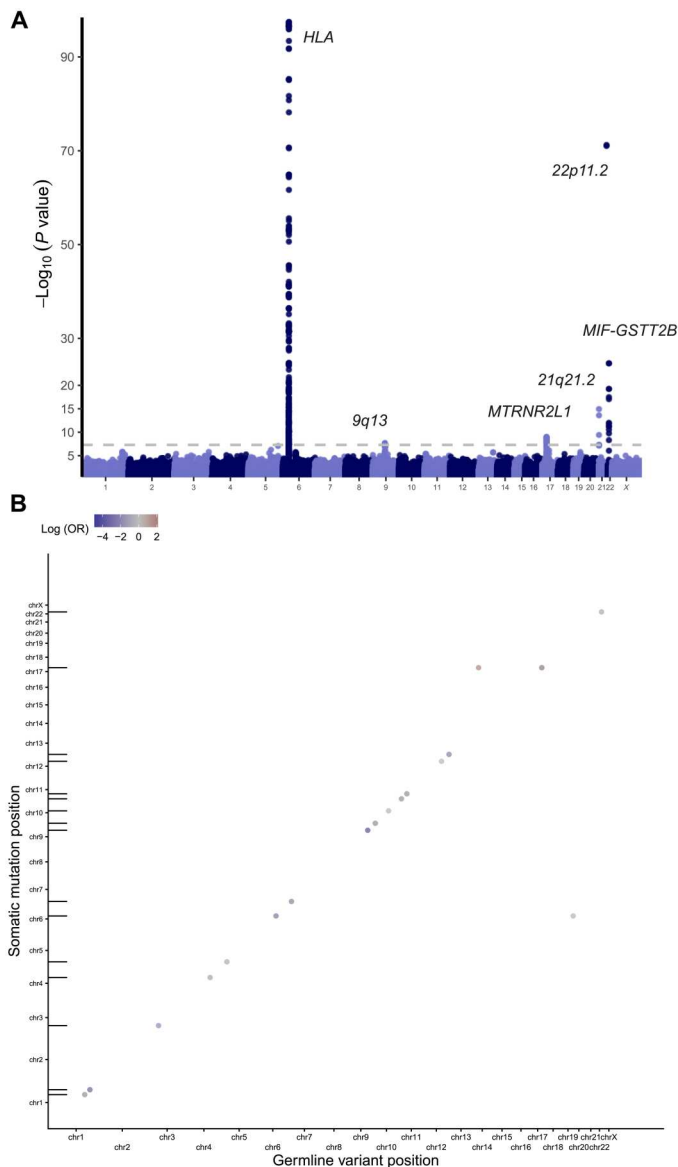


Fig. 3. Germline determinants of RNMSM burden. (A) Manhattan plot from GWAS of RNMSM burden, computed using SAIGE. Germline variants included had a minor allele count ≥ 600 and were distinct from the set of RNMSMs. (B) Genetic determinants of individual RNMSMs are primarily in “cis,” where cis is defined as within 2 Mb of the RNMSM. Genomic coordinates of the associated linkage disequilibrium (LD)-clumped germline variants are plotted on the x axis, and coordinates of the RNMSMs are plotted on the y axis. Points falling along the diagonal indicate cis associations. Tick marks along the y axis indicate positions of the RNMSMs from the LD-clumped associations.

possibly overlapping windows for association (82 2-Mb regions \times 82 RNMSMs) and used 382,359 germline variants (31,353,438 total germline variant–RNMSM pairs).

We identified 3094 associations at P value $< 5 \times 10^{-9}$, comprising 2783 germline variants and 79 loci (tables S12 and S13). Consistent with prior observations regarding somatic structural variants, such as mosaic chromosomal abnormalities (mCAs) (44), the vast majority (99.8% or 3088 of 3094) of significant hits were in cis (Fig. 3B).

Nine of 2783 germline variants altered protein coding sequences, including eight missense variants and one splice donor variant. The C allele of rs3737737, a missense variant in *MAST2*, strongly reduced the odds of acquiring a mutation in cis [chr1-46207599-T-C, odds ratio (OR) = 0.26, P value = 8.7×10^{-12}]. *MAST2* is a serine/threonine kinase that regulates interleukin-12 synthesis in macrophages (46). The G allele of rs3737744, a missense variant in *NSUN4*, was associated with increased odds of acquiring the same mutation (chr1-46207599-T-C, OR = 1.27, P value = 8.4×10^{-11}). Both rs3737744 and rs3737737 have been previously associated with altered hematocrit levels in GWAS of blood cell counts (37).

We then examined noncoding associations and identified four variants (rs116227346, rs1485836, rs7152548, and rs113359887) with a scaled combined annotation dependent depletion (CADD) (47) score ≥ 20 , indicating that they are predicted to be among the 1% most deleterious mutations. According to chromHMM (48) annotations, rs116227346 is in an enhancer in 24 epigenomes, including CD14-positive monocytes and common myeloid progenitors. rs1485836 lies in a region of heterochromatin in a highly conserved sequence (PhyloP = 2.5) 173 kb from the TSS of *GRIK2*. Among other hits, the C allele of rs1707302 was also associated with reduced odds of mutated chr1-46207599-T-C (OR = 0.71, P value = 3.6×10^{-18}). rs1707302 is 2.2 kb from the TSS of *PIK3R3* and is an eQTL in monocytes of *PIK3R3* (49) and also an sQTL of *MUTYH* in whole blood (35). *MUTYH*, which is involved in base excision repair, has recently been shown to alter somatic mutation rates in normal human cells (50) and alter the germline mutation rate in mice (51). The C allele of rs116750427, an intronic variant of *PAX5*, was associated with reduced odds of acquiring a mutated chr9-37358292-A-G (OR = 0.61, P value = 5.7×10^{-10}). The C allele is common in African ancestry genomes [3.5% in 1000 Genomes (52)] but unobserved in other ancestries. *PAX5* is a transcription factor that regulates lineage commitment in lymphoid progenitors and regulates VDJ recombination (53). The G allele of rs112886224, a variant upstream of *MDM2*, was associated with reduced odds of acquiring a mutated chr12-68547943-G-T. *MDM2* is a proto-oncogene that negatively regulates p53 (54). Together, these genetic association analyses reveal that inherited genetic variation shapes somatic mutations through three distinct mechanisms. First, inherited genetic mutations can result in an increased likelihood of somatic mutations nearby through a hypermutable local milieu. Second, germline genetic variants can lead to alterations in adaptive immune cell function, as suggested by the HLA signal. Third, germline genetic variants can promote an inflammatory state, which may lead to accelerated aging of HSCs, as suggested by the *MIF-GSTT2B* signal.

Phenotypic consequences of RNMSMs

We then sought to establish whether RNMSMs have any human health consequences. We developed an RNMSM association study (RSMAS) procedure, where we examined the association between each somatic variant with a given phenotype. To mitigate confounding, we included an extensive covariate adjustment set, including separate PCs from both the germline genetic variation (gPCs) and the somatic variation (sPCs), as well as measures of sample sequencing depth and contamination estimates (Methods). As previous reports have implicated mosaicism of driver mutations with altered hematopoiesis (8), and our germline

genetic studies highlighted genes important in blood cell function, we focused our RSMAS on 15 blood cell traits (Methods).

We performed 106,955 trait-variant associations between 15 blood cell traits and the 7131 RNMSMs. To adjust for multiple testing, we first identified the effective number independent RNMSMs using the eigenvalue ratio measure (55). The effective number of independent tests (4598) is lower than the number of tested RNMSMs (7131) because of the correlation structure present, resulting in a significance threshold of 1.1×10^{-5} for each individual blood trait. To estimate the false discovery rate of these associations, we also computed q values (56). Eight significant associations were detected at a q value $< 5\%$ threshold, and four were detected with a P value $< 1.1 \times 10^{-5}$ (Fig. 4). Of the eight q value significant variant-trait pairs, six were RNMSM-monocyte count associations, and one association was observed for hemoglobin and hematocrit, respectively. The eight RNMSMs were all rare, each with fewer than 100 mutations. The effect size of these associations is large compared to common germline variant disease associations: An increase in VAF of 0.1 in rs1319097708 is associated with an increase in monocyte count by half an SD [$\beta = 0.55$, 95% CI (0.301, 0.804), q value = 0.036] and explains 0.21% of the variance in monocyte counts.

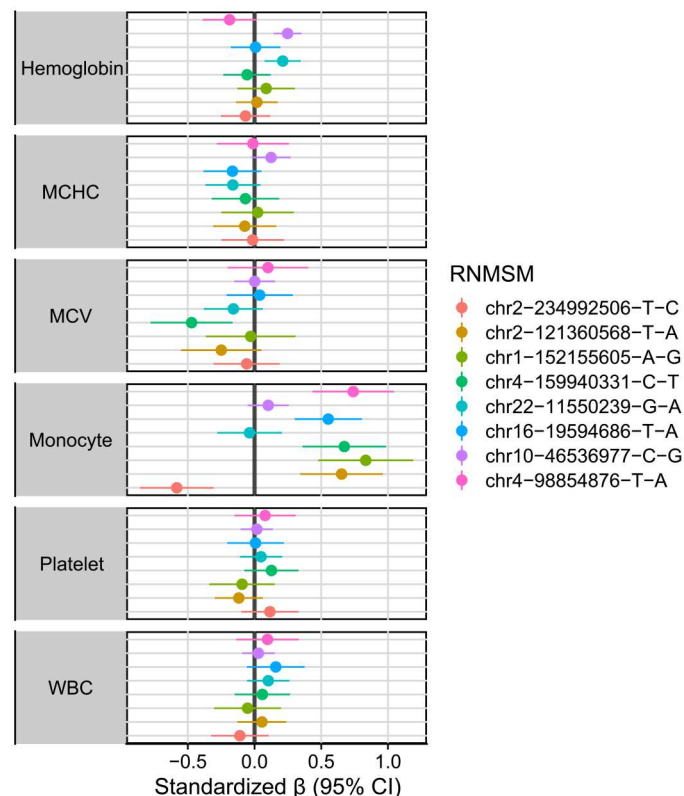


Fig. 4. RNMSMs associate with blood cell traits. Associations between eight RNMSMs (q values < 0.05) with harmonized blood cell traits in 2996 individuals. Units correspond to an increase of an SD of the phenotype associated with an increase in the VAF value of an RNMSM by 0.10. WBC, white blood cell; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume.

DISCUSSION

In this study, we characterized the spectrum of somatic variation in blood in 43,693 WGS individuals from 37 studies that comprise the TOPMed WGS dataset. We observed that most somatic variants are individually very rare. However, there were 7131 non-missense mutations present in more than 50 genomes, corresponding to an MAF of $>0.05\%$. RNMSM burden increased with age. These variants were highly enriched in heterochromatin. A notable number of these were differentially present in individuals of different genetic ancestries, due in part to ancestry-specific cis and trans QTL variants that associate with RNMSMs. The RNMSM-QTLs implicated diverse mechanisms including lymphoid lineage commitment, inflammatory cytokine production, and DNA repair. Eight RNMSMs were significantly associated with blood cell traits.

These findings permit several conclusions. First, nononcogenic somatic variation is a common phenomenon. Although the majority of somatic variants are individually very rare, a subset are quite common. Every sample had at least one RNMSM mutation. Analysis of other somatic mutations in blood cells, such as clonal hematopoiesis, suggests that they reoccur as a consequence of mutation rate and fitness (57). Given the limited functional consequence of RNMSMs, we identify three possible causes of this phenomenon. First, they reoccur at sites with high mutation rates in hematopoietic cells. This suggests an etiology where multiple HSCs acquire the same mutation rather than a single clone acquiring an RNMSM and then clonally expanding. Second, a subset of RNMSMs occur in B cell or T cell progenitors that subsequently expand clonally. Third, a subset of RNMSMs may cooperate with driver mutations to promote clonal expansion.

Second, the spectrum of somatic variation in humans is associated with genetic ancestry. Concordantly, two recent reports have indicated that germline variation is associated with the acquisition of somatic copy number variants in both the U.K. Biobank (45) and BioBank Japan (32), and they observed considerable differences in the prevalence of specific mCAs. Our observation extending this to somatic point mutations has implications for our understanding of cancer biology and its treatment. Recent reports have indicated that tumor mutational burden, a biomarker of response to cancer immunotherapy, is associated with germline variation (58). Differences in predisposition toward mutational burden is derived in part from differences in germline variation mediating somatic variation, highlighting a potential disparity in immunotherapy outcomes and an important area for future research.

Third, noncoding somatic variants associate with blood cell traits. We also introduce RSMAS, an analogous analysis to a GWAS on germline variants. GWAS has underscored the importance of regulatory biology in complex trait variation. In contrast to germline association studies, little is known about regulatory somatic variation. Despite the enrichment of germline loci in regulatory regions, most somatic analyses have focused on the acquisition and identification of oncogenic drivers rather than sets of variants with modest effects that may collectively contribute non-trivial variation to phenotypes. We anticipate that this method will yield novel insights in other phenotypes and tissues.

This study has several limitations. First, our variant-calling procedures are limited by the sensitivity of the sequencing used. Some sequencing artifacts have similar properties to somatic mutations, and a subset of somatic mutations are difficult to distinguish

from sequencing and mapping artifacts. We observed that our variant calling required extensive filtering to exclude likely artifacts. Refined somatic variant calling processes for single-tissue cohorts remain an area of future research. At 38× coverage, we are also insensitive to low VAF (<5%) variants (8). However, previous reports suggest that these presently undetected variants have little bearing on disease (59). This suggests that variants present in our RNMSMs are likely enriched for disease associations relative to the undetected class of low VAF RNMSMs. Second, as our compendium is derived from a complex collection of cohorts, ethnicities, and disease sampling ascertainment, our results are limited by our relatively incomplete knowledge of potential environmental confounders. Third, somatic variants are more challenging to interpret than germline variants, as causal inference is difficult to achieve in cross-sectional analyses, where date of mutation acquisition is unknown. In the absence of longitudinal or other study designs that permit causal inference in this setting, we cannot fully characterize the direction of causality in associations between RNMSMs and phenotypes. Fourth, as few other population cohorts with whole-genome somatic variant calls exist, it is difficult to replicate our observations in other contexts. Although most RNMSMs were found in gnomAD, a subset of the most frequent RNMSMs were not observed in a small, paired heart-blood tissue cohort. Further research is needed to fully characterize the technical factors that contribute to differences between paired-tissue and single-tissue somatic variant calling. Fifth, because of the absence of somatic variants in catalogs of molecular QTLs, the functional consequence of the blood-associated RNMSMs is unclear. Sixth, although we observed only a modest association with CHIP, it is possible that other forms of clonal hematopoiesis that we were unable to detect at our sequencing depth may contribute to the genesis of RNMSMs.

In conclusion, our compendium of recurrently mutated somatic variants demonstrates that this class of genetic variation is widely present throughout the human genome and has germline genetic determinants and phenotypic consequences. Our observations expand at scale our understanding of the spectrum of acquired human genetic variation and its origin. Future efforts are required to extend these observations to other human tissues.

METHODS

WGS processing, variant calling, and CHIP annotation

BAM files were remapped and harmonized through the functionally equivalent pipeline (60). SNPs and indels were discovered across TOPMed and were jointly genotyped across samples using the GotCloud pipeline (13, 61). A support vector machine (SVM) filter was trained to discriminate between high- and low-quality variants. Variants were annotated with snpEff 4.3 (62). Sample quality was assessed through Mendelian discordance, contamination estimates, sequencing converge, and among other quality control metrics. Samples were aligned to GRCh38.

Putative somatic SNPs were called with GATK Mutect2 (11), which searches for sites with at least two alt-reads present. Mutect2 then performs a local haplotype reassembly and applies several variant quality filters. We used a “panel of normals” to filter sequencing artifacts and used an external reference of germline variants to exclude germline calls. We deployed this pipeline on Google Cloud using Cromwell (63).

Samples were annotated as CHIP carriers if the Mutect2 output contained at least one variant in a curated list of leukemogenic driver mutations with at least three alt-reads supporting the call (7, 8). We expanded the list of driver mutations to include those in recently identified CHIP genes (64).

We defined RNMSMs as those occurring in at least 50 distinct samples and occurring in at least five TOPMed cohorts, which would indicate a prevalence similar to those of the most frequently mutated CHIP genes. We used *cyvcf2* (65) to parse the Mutect2 VCFs and encoded each variant in an int64 value using the variant key encoding (66). We developed a Python application to perform the recurrent somatic variant identification and filtering. We stored the VAFs values in a sparse matrix using the Python Sparse module. For a given individual and sample, we stored the VAF of the mutation if present and, otherwise, coded the value as 0. We performed extensive filtering to exclude germline variants and sequencing artifacts (text S1).

RNMSM validation in paired tissue samples

We gathered paired blood and heart tissue samples from 46 donors from Vanderbilt University Medical Center. We sequenced the whole genomes of the 92 samples at ~30× coverage and called somatic variants using the Illumina DRAGEN small variant caller in paired tumor/normal mode. We examined variants that were marked as either “PASS” or “weak evidence” as indicating replication for the RNMSMs of our discovery cohort. We limited our analysis to only those variants that were present in the blood tissues with at least one minor allele. We then examined the intersection between the RNMSM site list from our discovery cohort with the validation samples using *bctools* (67).

VEP annotation and enrichment

We annotated the RNMSMs using VEP (16) using the “regulatory” and CADD plugins. We filtered the annotations using the “flag_pick” command line option. As a comparison set, we applied the same VEP command to all PASS variants from chromosome 22 of the TOPMed Freeze 8 binary variant call format (BCF) files. We then compared the proportions of variants with a given annotation using a proportion *Z* test.

Roadmap Epigenomics enrichment

We downloaded the hg38 mnemonic files from the Roadmap Epigenomics project derived from the chromHMM 15-state model applied to 127 epigenomes (48, 68). We then constructed a searchable index over the 1905 (51 * 127) annotations using GIGGLE (69). We used GIGGLE to calculate the genome-wide enrichments of the RNMSMs across the 1905 annotations. GIGGLE implements the enrichment hypothesis testing using a Fisher’s exact test.

Mutational signature analysis

We estimated the contributions of the v3 COSMIC mutational signatures to the RNMSMs using SignatureAnalyzer (70). We set “nruns” = 50 and the “a” hyperparameter to 15.

PCA on RNMSMs

We computed the partial singular value decomposition (SVD) of the sparse matrix of VAFs using the augmented implicitly restarted Lanczos bidiagonalization algorithm as implemented in the IRLBA R package (26). We normalized the VAF matrix by column

centering and scaling by the VAF SD for the given RNMSM. We computed the first 50 nonzero components among the left and right singular vectors.

Fst calculation

We stratified samples into European, Sub-Saharan Africa, and East Asian ancestry groups based on their RFMix global ancestry estimate, which was defined as the largest global ancestry estimate for each sample. Within each strata, we computed the allele frequencies of all RNMSMs separately. We then computed the RNMSM variance within each strata k as $\text{var}_k = 2 * \text{AF}_k(1 - \text{AF}_k)$. We then computed the weighted average of the stratified variance estimates, $\overline{\text{var}} = \sum c_k \text{var}_k$, where c_k indicates the prevalence of group k . We calculated $F_{\text{st}} = \frac{\text{var}_{\text{tot}} - \overline{\text{var}}}{\text{var}_{\text{tot}}}$, where var_{tot} is the binomial variance of the RNMSM computed across all samples.

Single variant association analyses on RNMSM burden

Single variant association for each variant in Freeze 8 with a minor allele count (MAC) > 20 was performed with SAIGE (33) using the TOPMed Encore analysis server. We defined the quantitative RNMSM burden phenotype by counting the RNMSM mutations per sample and then applying an inverse normal transformation. We included age at blood draw, genotype-inferred sex, study, average sample depth, sample contamination from VerifyBamID (71), sPCs 4-5, and the first 10 genetic ancestry PCs as covariates. We declared variants from this analysis as significant if their P value was less than 5×10^{-8} .

We also computed European ancestry- and Sub-Saharan ancestry-specific GWAS by first stratifying the samples based on their RFMix global ancestry label. We then repeated the same GWAS analysis within each strata separately.

As the TOPMed germline variant call set intersects with the somatic calls, we took additional filtering measures to exclude somatic variants from the germline variant calls. First, we excluded any variant that was also called as an RNMSM. Second, we then created a database of variants where either the allelic-balance z score exceeded 5.0, the SVM filter was less than -0.25 , or the $-\log_{10}$ Hardy-Weinberg equilibrium (HWE) P value was greater than 5.0. We then excluded any variant in the summary statistics where at least one of these criteria was true. Last, we included only variants where the minor allele count was at least 600.

Single variant association analyses on individual RNMSMs

We selected the 80 most frequent RNMSMs and performed association analyses between each germline variant in a 2-Mb region surrounding each RNMSM using PLINK2 (72) with the “first-fallback” option, which defaults to standard logistic regression to compute statistics but switches to firch logistic regression when convergence fails. To identify trans RNMSM regulators, we also tested each of the RNMSM cis regions with all the other RNMSMs. In total, we examined 6724 possibly overlapping windows for association (82 2-Mb regions \times 82 RNMSMs). We assessed significance at 5×10^{-9} . As with the analysis of RNMSM burden, we then filtered the summary statistics to exclude somatic variants from the germline variant calls and to exclude any variant with a MAC less than 700.

Association analyses with RNMSMs (RSMAS)

For quantitative outcomes, we computed test statistics using linear regression on the inverse normal transformed outcome. We computed a marginal analysis where each RNMSM was analyzed in a separate regression. Each RNMSM was included as a vector of VAF values. We included the first 10 germline ancestry PCs, the first 10 sPCs, an indicator for study, genotype-inferred sex, average sample depth, and average sample contamination as covariates. We calculated q values (56) on the resultant test statistics within each outcome. We implemented this procedure in a bespoke R package, somaticWAS.

Blood cell traits

Blood cell counts and indices were selected for analysis including hemoglobin, hematocrit, red blood cell count, white blood cell count, basophil count, eosinophil count, neutrophil count, lymphocyte count, monocyte count, platelet count, mean corpuscular hemoglobin, mean corpuscular hemoglobin concentration, mean corpuscular volume, mean platelet volume, and red cell distribution width. The phenotypes were collected by each cohort and were centrally harmonized by the TOPMed Data Coordinating Center. Before analysis, each blood cell trait was separately inverse normal transformed.

Code availability

Code availability can be found at Zenodo archives: 10.5281/zenodo.7484723 and 10.5281/zenodo.7484713

Supplementary Materials

This PDF file includes:

Figs. S1 to S3

Variant Calling Supplementary Text

TOPMed Cohort Acknowledgements

Legends for tables S1 to S13

Other Supplementary Material for this manuscript includes the following:

Table S1 to S13

[View/request a protocol for this paper from Bio-protocol.](#)

REFERENCES AND NOTES

1. K. Yizhak, F. Aguet, J. Kim, J. M. Hess, K. Kübler, J. Grimsby, R. Frazer, H. Zhang, N. J. Haradhvala, D. Rosebrock, D. Livitz, X. Li, E. Arich-Landkof, N. Shores, C. Stewart, A. V. Segrè, P. A. Branton, P. Polak, K. G. Ardlie, G. Getz, RNA sequence analysis reveals macroscopic somatic clonal expansion across normal tissues. *Science* **364**, eaaw0726 (2019).
2. P. Desai, N. Mencia-Trinchant, O. Savenkov, M. S. Simon, G. Cheang, S. Lee, M. Samuel, E. K. Ritchie, M. L. Guzman, K. V. Ballman, G. J. Roboz, D. C. Hassane, Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nat. Med.* **24**, 1015–1023 (2018).
3. S. Abelson, G. Collord, S. W. K. Ng, O. Weissbrod, N. Mendelson Cohen, E. Niemeyer, N. Barda, P. C. Zuzarte, L. Heisler, Y. Sundaravadanam, R. Luben, S. Hayat, T. T. Wang, Z. Zhao, I. Cirlan, T. J. Pugh, D. Soave, K. Ng, C. Latimer, C. Hardy, K. Raine, D. Jones, D. Hoult, A. Britten, J. D. McPherson, M. Johansson, F. Mbabaali, J. Eagles, J. K. Miller, D. Pasternack, L. Timms, P. Krzyzanowski, P. Awadalla, R. Costa, E. Segal, S. V. Bratman, P. Beer, S. Behjati, I. Martincorena, J. C. Y. Wang, K. M. Bowles, J. R. Quirós, A. Karakatsani, C. La Vecchia, A. Trichopoulou, E. Salamanca-Fernández, J. M. Huerta, A. Barricarte, R. C. Travis, R. Tumino, G. Masala, H. Boeing, S. Panico, R. Kaaks, A. Krämer, S. Sieri, E. Riboli, P. Vineis, M. Foll, J. McKay, S. Polidoro, N. Sala, K.-T. Khaw, R. Vermeulen, P. J. Campbell, E. Papaemmanuil, M. D. Minden, A. Tanay, R. D. Balicer, N. J. Wareham, M. Gerstung, J. E. Dick, P. Brennan, G. S. Vassiliou, L. I. Shlush, Prediction of acute myeloid leukaemia risk in healthy individuals. *Nature* **559**, 400–404 (2018).

