

# UCSF

## UC San Francisco Previously Published Works

### Title

Inherited causes of clonal haematopoiesis in 97,691 whole genomes

### Permalink

<https://escholarship.org/uc/item/7w16b34n>

### Journal

Nature, 586(7831)

### ISSN

0028-0836

### Authors

Bick, Alexander G  
Weinstock, Joshua S  
Nandakumar, Satish K  
[et al.](#)

### Publication Date

2020-10-29

### DOI

10.1038/s41586-020-2819-2

Peer reviewed



Published in final edited form as:

Nature. 2020 October ; 586(7831): 763–768. doi:10.1038/s41586-020-2819-2.

## Inherited Causes of Clonal Hematopoiesis in 97,691 TOPMed Whole Genomes

A full list of authors and affiliations appears at the end of the article.

### Abstract

Age is the dominant risk factor for most chronic human diseases; yet the mechanisms by which aging confers this risk are largely unknown.<sup>1</sup> Recently, the age-related acquisition of somatic mutations in regenerating hematopoietic stem cell populations leading to clonal expansion was associated with both hematologic cancer<sup>2–4</sup> and coronary heart disease<sup>5</sup>, a phenomenon termed ‘Clonal Hematopoiesis of Indeterminate Potential’ (CHIP).<sup>6</sup> Simultaneous germline and somatic whole genome sequence analysis now provides the opportunity to identify root causes of CHIP. Here, we analyze high-coverage whole genome sequences from 97,691 participants of diverse ancestries in the NHLBI TOPMed program and identify 4,229 individuals with CHIP. We identify associations with blood cell, lipid, and inflammatory traits specific to different CHIP genes. Association of a genome-wide set of germline genetic variants identified three genetic loci associated with CHIP status, including one locus at *TET2* that was African ancestry specific. *In silico*-informed *in vitro* evaluation of the *TET2* germline locus identified a causal variant that disrupts a *TET2* distal enhancer resulting in increased hematopoietic stem cell self-renewal. Overall, we observe that germline genetic variation shapes hematopoietic stem cell function

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Please address correspondence to: Pradeep Natarajan, MD MMSc, 185 Cambridge St, CPZN 3.184, Boston, MA 02114 USA, [pradeep@broadinstitute.org](mailto:pradeep@broadinstitute.org), Twitter: @pnatarajanmd, Sekar Kathiresan, MD, 75 Ames St, Cambridge, MA 02139 USA, [sekar@broadinstitute.org](mailto:sekar@broadinstitute.org), Twitter: @skathire.

‡Present Address: Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, TN.

#### AUTHOR CONTRIBUTIONS

A.G.B., P.N. and S.K. conceived the study. A.G.B. and J.S.W. performed the germline and somatic whole genome sequence analyses. C.P.F., E.L.B., S.M.Z., M.D.S., M.J.L., J.N., K.C., C.J.G., A.E.L., B.B.B., P.S., J.K., J.M.E., A.P.R., B.L.E., and S.J. performed additional bioinformatic analyses. S.K.N., X.L. and V.G.S. experimentally characterized the *TET2* locus. M.A.T., F.A., K.A., B.D.M., K.C.B., A.M., M.F., S.R., B.M.P., E.K.S., S.T.W., N.D.P., R.S.V., E.G.B., S.L.R.K., J.H., R.C.K., N.L.S., D.K.A., D.A.S., A.C., M.d.A., X.G., B.A.K., B.C., J.M.P., H.G., D.A.M., S.T.M., I.Y., M.B.S., P.A.P., J.G.B., S.M.G., F.F.W., Q.W., M.E.M., M.D., E.E.K., K.E.N., L.J.L., B.E.C., J.C.B., M.H.C., J.L.S., D.W.B., L.A.C., A.C.M., L.C.B., J.A.S., T.N.K., S.A., S.R.H., H.K.T., I.V.Y., J.A.H., S.L., J.M.J., J.E.C., S.E.W., D.E.W., D.C.R., D.D., J.Y.M., R.P.T., E.J.B., N.R., R.J.F.L., P.D., Y.L., L.H., J.L., P.K., B.I.F., D.L., L.F.B., J.E.H., J.S.F., E.A.W., P.T.E., M.R.I., T.E.F., L.M.R., S.M.A., M.M.W., E.C.S., J.B., L.K.W., B.D.L., W.H.S., D.M.R., E.B., J.E.M., R.A.M., P.D., K.D.T., A.D.J., P.L.A., C.K., C.C.L., T.W.B., A.V.S., H.Z., E.L., L.L., S.S.R., J.I.R., J.G.W., P.S., J.O.K., E.S.L., J.M.E., and G.A. contributed to sample acquisition, DNA sequencing and phenotypic curation for the NHLBI TOPMed constituent cohorts analyzed here. A.G.B., J.S.W., S.K. and P.N. wrote the manuscript with input from all authors.

\*These individuals contributed equally to this work

#These individuals jointly supervised this work

†A list of authors and their affiliations appears at the end of the paper

#### DATA AVAILABILITY

Individual whole-genome sequence data for TOPMed whole genomes, individual-level harmonized phenotypes, harmonized germline variant call sets, the CHIP somatic variant call sets, RNA-Seq and peripheral blood methylation data used in this analysis are available through restricted access via the dbGaP. Accession numbers for these datasets are provided in Supplementary Table 1. Summary-level genotype data are available through the BRAVO browser (<https://bravo.sph.umich.edu/>). Full GWAS summary statistics are available for general research use through controlled access at dbGaP accession phs001974: NHLBI TOPMed: Genomic Summary Results for the Trans-Omics for Precision Medicine Program. A subset of the TOPMed cohorts analyzed here are based on sensitive populations, precluding public sharing of full genomic summary results.

leading to CHIP through mechanisms that are both specific to clonal hematopoiesis and shared mechanisms leading to somatic mutations across tissues.

---

The U.S. National Heart, Lung, and Blood Institute (NHLBI) Trans-omics for Precision Medicine (TOPMed) project seeks to use high-coverage (>35x) whole genome sequencing (WGS) and molecular profiling to improve fundamental understanding of heart, lung, blood, and sleep disorders.<sup>7</sup> Within the TOPMed program, we designed a study to detect CHIP from blood DNA-derived WGS in 97,691 individuals across 52 largely observational epidemiologic studies to discover the inherited genetic causes and phenotypic consequences of CHIP (Supplementary Table 1).

To confidently identify somatic mutations in blood-derived DNA, we applied a somatic variant caller<sup>8</sup> to TOPMed WGS data. We identified CHIP carriers on the basis of a pre-specified list of leukemogenic driver mutations (see Methods, Supplementary Table 2).<sup>5</sup>

In total, we identified 4,938 CHIP mutations in 4,229 individuals (Supplementary Table 3). The median variant allele fraction (VAF) of the CHIP mutations observed was 16%. Consistent with prior reports, >75% of these CHIP mutations were in one of three genes, *DNMT3A*, *TET2*, and *ASXL1*. Approximately 15% of these CHIP mutations were in the five next most frequent genes (*PPM1D*, *JAK2*, *SF3B1*, *SRSF2* and *TP53*, Figure 1). Amongst these 8 genes, there was marked heterogeneity in clonal fraction. For example, *DNMT3A* and *TET2* CHIP clonal fraction of the peripheral blood was ~25% smaller ( $p=1.3 \times 10^{-15}$ ) and ~14% smaller ( $p=2.1 \times 10^{-4}$ ), respectively, than *ASXL1* clonal fraction, implicating the presence of driver mutation gene-specific differences in clonal selection (Extended Data Figure 1a). 90% of individuals with CHIP driver mutations had only one identified mutation (Extended Data Figure 1b).

## CHIP phenotypic associations

CHIP prevalence was strongly correlated with age at blood draw ( $p < 10^{-300}$ , Figure 1 inset). CHIP prevalence was highly consistent across studies and comparable to previous reports<sup>2-4</sup> using whole exome sequencing (Extended Data Figure 1c,d). Consistent with prior studies, history of smoking was associated with increased CHIP odds (OR = 1.18,  $p=5 \times 10^{-5}$ ) whereas Hispanic ancestry and East Asian ancestry were each associated with reduced CHIP odds (OR = 0.50,  $p=0.008$  and OR = 0.56,  $p=0.001$  respectively) after adjusting for age (Supplementary Table 4).

Carriers of frameshift CHIP mutations were on average older individuals than carriers of single nucleotide CHIP mutations (Wilcoxon rank sum test:  $p=0.01$ ). In the subset of individuals with *ASXL1* CHIP mutations, which are exclusively loss-of-function single nucleotide stop-gain mutations or frameshift mutations, *ASXL1* frameshift mutation carriers were similarly older (Wilcoxon rank sum test:  $p=0.009$ , Extended Data Figure 3a).

*JAK2* CHIP carriers were the youngest among CHIP carriers. Relative to *JAK2*, *ASXL1* and *TET2* carriers were 3.3 ( $p=0.01$ ) and 3.9 ( $p=9.1 \times 10^{-4}$ ) years older, respectively, while

*PPM1D*, *SF3B1* and *SRSF2* carriers were 5.0, 6.9 and 7.7 years older ( $p=5.7 \times 10^{-4}$ ,  $1.8 \times 10^{-6}$ ,  $1.3 \times 10^{-4}$ ), respectively (Extended Data Figure 3b).

To evaluate the overlap between CHIP and large-scale mosaic chromosomal rearrangements<sup>9</sup>, we evaluated a subset of 855 samples with both WGS and array genotyping data. The two somatic events did not co-occur more than expected by chance (hypergeometric  $p=0.25$ , Extended Data Figure 3c).

CHIP is distinguished from other clonal hematologic disorders based on the absence of cytopenia, dysplasia, and neoplasia.<sup>6</sup> We observed a modest increase in total white blood cell count ( $p=1.1 \times 10^{-5}$ ) and a modest decrease in hemoglobin ( $p=0.04$ ), among those with CHIP compared to those without (Extended Data Figure 3a, Supplementary Table 5). In aggregate, CHIP driver mutations were associated with increased red blood cell distribution width (RDW,  $p=3.0 \times 10^{-5}$ ) consistent with prior observations.<sup>10</sup> Notably, RDW is a hematologic parameter that increases with age and predicts overall mortality and poor clinical outcomes in the setting of CVD and in older adults.<sup>11</sup>

Given the prior association of CHIP with atherosclerotic cardiovascular disease<sup>5,12</sup>, we asked whether CHIP carriers had altered lipid profiles. Consistent with prior reports<sup>5</sup>, we observed negative correlations of *JAK2* CHIP carrier status with total cholesterol ( $p=5.1 \times 10^{-4}$ ) and LDL cholesterol ( $p=0.0014$ ) but no other significant associations (Extended Data Figure 3b, Supplementary Table 6).

We characterized the human inflammatory profile of CHIP carriers (Extended Data Figure 3c, Supplementary Table 7). In aggregate, CHIP was associated with increased IL-6 ( $p=0.0035$ ). There was no association of CHIP with quantitative C-reactive protein (CRP) and elevated CRP did not reliably identify carriers of CHIP (AUC: 0.55; for cutoff of  $CRP > 2$  mg/L: PPV=6.3%, sensitivity=60%). Driver gene-specific analyses highlighted the association of *TET2* CHIP with increased IL-1b ( $p=2.4 \times 10^{-4}$ ), while *JAK2* and *SF3B1* were associated with increased circulating IL-18 ( $p=1.3 \times 10^{-4}$  and  $1.27 \times 10^{-20}$  respectively).

To identify underlying determinants of the somatic mutational spectrum, we performed COSMIC mutational signature analysis<sup>13</sup> on passenger somatic mutations identified in CHIP carriers and non-carriers (see Methods). Among CHIP carriers, we observe enrichment of signature 4, which has been associated with smoking, as well as signature 6, which has been associated with defective DNA mismatch repair. (Extended Data Figure 5).

## Germline genetic determinants of CHIP

Germline genetic variants have been previously associated with clonal hematopoiesis, defined either by somatic mosaicism of SNVs and indels<sup>14</sup> or by large scale chromosomal rearrangements<sup>9</sup>, in individuals of European ancestry, and identified variants at a single locus, *TERT*, that associates with clonal hematopoiesis. Given the distinct association of clonal hematopoiesis with known leukemogenic mutations (i.e., CHIP) with both cancer<sup>2,15,16</sup> and atherosclerotic cardiovascular disease<sup>5,12</sup>, we sought to discover germline genetic variations conferring increased risk for CHIP acquisition. We performed a single variant genome-wide association analysis in a subset of 65,405 individuals (3,831 CHIP

cases) where the likelihood of having a CHIP mutation was >1% (see Methods). The trait heritability explained by the analysis with LD score-regression was 3.6%.

Our WGS-based association analysis of CHIP replicated the lead variant of the single locus previously associated at genome wide significance with clonal hematopoiesis (defined based on somatic mosaicism of SNVs and indels),<sup>14</sup> rs34002450 (OR 1.2,  $p=2.0 \times 10^{-13}$ ). rs34002450 is in strong LD ( $r^2=0.55$ ) with our lead variant at this locus rs7705526, a common variant (MAF 0.29) in the 5<sup>th</sup> intron of *TERT*, which encodes telomere enzyme reverse transcriptase. In TOPMed, carriers of the rs34002450-A (minor) allele have a 1.3-fold risk of developing CHIP ( $p=8.4 \times 10^{-24}$ ). This variant was previously significantly associated with increased leukocyte telomere length<sup>17</sup>, myeloproliferative neoplasms (MPN, Bao, co-submitted manuscript) and clonal chromosomal mosaicism<sup>9</sup>. In a phenome-wide association analysis (PheWAS) of rs34002450-A in UK Biobank, we identified significant increased risk of MPN ( $p=2.6 \times 10^{-13}$ ), uterine leiomyoma ( $p=3.2 \times 10^{-9}$ ) and brain cancer ( $p=3.6 \times 10^{-8}$ ).

We performed a conditional analysis at the *TERT* locus, and identified a second intronic *TERT* variant rs13167280 (MAF 0.11,  $r^2=0.2$  with rs7705526) that independently associates with CHIP status (OR 1.3,  $p=6.1 \times 10^{-10}$ ; conditional OR: 1.1,  $p=4.7 \times 10^{-4}$ ).

In the TOPMed single-variant association analysis, we additionally identified 2 other novel genome-wide significant genetic loci, including one locus on chromosome 3 in an intergenic region spanning *KPNA4/TRIM59* and one locus on chromosome 4 near *TET2* (Figure 3, Extended Data Figure 6, Supplementary Table 8).

rs1210060191 is a common variant (MAF 0.54) in a locus with an association signal that spans a 300kb region that includes *KPNA4*, *TRIM59*, *IFT80*, and *SMC4*. The lead variant is a 1 bp intronic deletion in *TRIM59*. Carriers of the del(T) allele have a 1.16-fold increased risk of CHIP ( $p=5.3 \times 10^{-10}$ ). Variants in LD with this variant have been identified as associated with MPN (Bao et al, co-submitted manuscript). No other significant phenotypic associations were noted in UK Biobank PheWAS analyses.

rs144418061 is an African ancestry specific variant (MAF 0.035 in African Ancestry samples, not present in non-African-ancestry samples) in an intergenic region near *TET2*. Carriers of the A allele have a 2.4-fold increased risk for CHIP ( $p=4.0 \times 10^{-9}$ ). We replicated this association in an additional set of 570 TOPMed CHIP cases and 8,819 TOPMed controls (OR: 2.1,  $p=0.026$ ). The association is equally robust for *DNMT3A* CHIP, *TET2* CHIP and *ASXL1* CHIP, suggesting that the germline variant does not specifically predispose to *TET2* CHIP. Although other variants in the vicinity of *TET2* have been associated with MPN (Bao et al, co-submitted manuscript), this variant has not been previously identified as associated with any traits in the literature likely due to the under-representation of African ancestry genomes in published association studies.

We considered whether there might be germline variants that predispose to specific CHIP driver mutations by separately performing a GWAS on *DNMT3A* CHIP and *TET2* CHIP. We identified a single novel locus for *DNMT3A* chip at rs2887399 in an intron of T-cell leukemia/lymphoma 1A (*TCL1A*). Carriers of the T allele (MAF 0.26) are at 1.23-fold risk

of acquiring a *DNMT3A* CHIP mutation ( $p=3.9 \times 10^{-9}$ ). Intriguingly carriers of the T allele are at decreased risk of acquiring a *TET2* CHIP mutation (OR: 0.82,  $p=0.0012$ ), and consequently it was not identified in the primary CHIP GWAS analysis. This variant is also associated with mosaic loss of chromosome Y.<sup>18</sup>

We evaluated whether our associations between germline loci and CHIP clones were robust across CHIP clone size spectrum, using the association between the *JAK2* 46/1 haplotype (tagged by rs1327494) and *JAK2* CHIP.<sup>19</sup> We find that rs1327494 associates with *JAK2* CHIP presence across VAF thresholds. We evaluated whether this observation generalized beyond *JAK2* CHIP to encompass all CHIP. Intriguingly, we find that the *TERT* locus (tagged by rs7705526) is associated with CHIP presence across all VAF thresholds (Supplementary Table 9). These observations imply that our genetic associations are not dependent on clone size detectable by deep-coverage whole genome sequencing.

As single-variant analyses have limited power to detect rare-variant associations, we next performed several types of variant aggregation association tests. First, we performed a transcriptome-wide association analysis to quantify the relationship between changes in gene expression and genetic predisposition to CHIP<sup>20</sup> (see Methods). This approach identified the Chr3 *KPNA4/TRIM59* locus and six additional loci including: *AHRR*, *ASL*, *KREM2*, *LEAP2*, *JSRP1*, *RASEF*. (Extended Data Fig. 7–8) *AHRR* directs hematopoietic progenitor cell expansion and differentiation.<sup>21</sup>

We also performed gene-based association tests for aggregations of rare (MAF<0.1%) putative loss-of-function (pLOF) germline variants within genes for CHIP presence. Although no genes reached exome-wide significance, the top associated gene was DNA damage repair gene *CHEK2* (OR 1.7,  $p=1.3 \times 10^{-5}$ , Supplementary Table 10). Rare germline variants in *CHEK2* are implicated in a diverse set of hematologic and solid tumor malignancies.<sup>22,23</sup> Common variants in *CHEK2* are associated with MPN<sup>19</sup> and a low frequency frameshift *CHEK2* is associated with somatic chromosomal mosaicism<sup>9</sup>. In recent experimental work, suppression of *CHEK2* in human cord blood Lin<sup>-</sup>CD34<sup>+</sup> cells increased cellular proliferation in long term culture. (Bao et al, co-submitted manuscript) These results suggest that while *CHEK2* while may ordinarily limit hematopoietic stem cell expansion, loss of *CHEK2* function may promote self-renewal increasing risk of CHIP.

We next sought to determine whether rare variants in non-coding regions associate with CHIP acquisition (see Methods). One set of variants in *HAPLN1* enhancers exceeded a p-value threshold of  $p<0.05$  after Bonferroni-correction (OR: 6.8,  $p=1.96 \times 10^{-5}$ , Supplementary Table 11). *HAPLN1* is an extracellular matrix protein, produced in bone marrow stromal cells that has previously been implicated in NF- $\kappa$ B signaling.<sup>24</sup>

We asked whether germline genetic variants might be associated with CHIP clonal expansion. No single variants or aggregated rare variants exceeded Bonferroni significance (Supplementary Table 12–13).

## TET2 CHIP risk locus characterization

Lastly, we bioinformatically and experimentally characterized the mechanism by which the non-coding African American-specific variant at the *TET2* locus influenced risk for CHIP. First, iterative conditional analyses at the locus suggested that there was most likely only a single causal variant. Fine-mapping prioritized 25 variants in the credible set (>99% posterior probability), none of which overlaps the coding sequence or promoter of a protein-coding gene.

We hypothesized that the causal variant affects an enhancer for *TET2* in hematopoietic stem cells, because heterozygous *Tet2* knockout in mice increases the self-renewal of hematopoietic stem cells *in vivo*<sup>25</sup> and recapitulates the clonal expansion observed in humans with somatic mutations in *TET2*.<sup>5,10</sup> Accordingly, we used the Activity-by-Contact (ABC) model to predict which noncoding elements act as enhancers in CD34+ hematopoietic stem and progenitor cells (HSPC, see Methods). Only a single variant (rs79901204) in this credible set overlapped an element predicted to regulate any gene, and that element was indeed predicted to regulate *TET2* expression. (Figure 3a, Supplementary Table 14) The T risk allele disrupts a consensus GATA/E-Box motif, likely resulting in reduced binding of the activating transcription factors GATA1 and GATA2 (Figure 3b,c).

We then evaluated whether rs79901204 affected *TET2* expression *in vivo* in human peripheral blood samples. We utilized whole blood RNAseq from 247 African American individuals, 16 of whom were heterozygotes for rs79901204 and one who was a homozygote. In these samples, the T risk allele led to a dose-dependent decrease in whole blood *TET2* expression (Beta: -0.27, SE: 0.11, two-sided linear mixed model p=0.012, Figure 3d). Therefore, we sought to test our hypothesis that that the rs79901204 risk allele acts to decrease the activity of this *TET2* enhancer and that decreased enhancer activity reduces expression of *TET2* *in vitro*.

To test whether rs79901204 affects enhancer activity, we tested a 600 base pair region containing the regulatory element using a plasmid-based luciferase enhancer assay in hematopoietic cells. The reference sequence activated luciferase expression by 118-fold (versus control constructs with no enhancer sequence), while the T risk allele activated expression by only 27-fold (two-sided t-test p=0.007, Figure 3e).

To test whether deletion of this enhancer would alter *TET2* gene expression, we performed deletion of the enhancer element in CD34+ HSPCs using a pair of CRISPR/Cas9 guides introduced as ribonucleoproteins, which resulted in decreased *TET2* expression after 48 hours (Figure 3f).

We then sought to establish the effect of decreased *TET2* expression on HSPC expansion using a colony forming unit cellular assay. Human HSPCs were electroporated with Cas9 targeting a coding region of *TET2* and *AAVS1* (a control locus) and plated for primary and secondary colony-forming assays (Figure 3g). To establish a dose response relationship, two *TET2* guides were used with differential editing efficiency (Figure 3h, Extended Data Figure 9). *TET2* coding disruption resulted in expanded secondary colony formation compared to *AAVS1* controls, with greater expansion identified in the *TET2* guide with greater editing



efficiency (Figure 3i). Thus, we demonstrate that reduction of *TET2* activity promotes self-renewal and proliferation of HSPCs, illustrating how both germline noncoding and somatic coding variation at this locus converge to affect *TET2* and influence the development of CHIP.

Given the established role of *TET2* in DNA de-methylation and our finding that rs79901204 is associated with decreased *TET2* expression (Figure 3d), we hypothesized that carriers of rs79901204 T allele may have altered peripheral blood methylation profiles. We performed a methylation-QTL analysis of 1747 African Americans and identified 597 genes across the genome with differentially methylated CpG loci associated with rs79901204 T carrier status. The most strongly differentially methylated sites were at the *TET2* locus itself. (Extended Data Fig. 10, Supplementary Table 15)

Our observations permit several conclusions. First, our sample size which is nearly an order of magnitude larger than prior CHIP analyses<sup>2,3,14</sup> enables refinement of CHIP phenotype associations at the level of CHIP driver genes. We find that considerable heterogeneity exists across CHIP phenotypes by driver gene. For example, IL-1b and IL-18 both activate through the inflammasome and increase IL-6. However, while *TET2* CHIP is associated with increased levels of IL-1b, *JAK2* and *SF3B1* CHIP are associated with IL-18.

Second, our work highlights multiple mechanisms through which germline genetic variation can shape somatic variation in hematopoietic stem cells. A set of the germline loci are associated with increased propensity to acquire mutations due to failure of genes that maintain genome integrity (e.g. *TERT* and *CHEK2*) and which have been implicated in stem cell maintenance/self-renewal (*Bao et al*, co-submitted manuscript). These loci are associated with acquisition of somatic mutations resulting in neoplasm in multiple tissues. Other germline loci are associated with increased hematopoietic stem cell self-renewal (e.g. *TET2*). While the *TET2* locus is associated with increased risk of acquiring any CHIP driver mutations, it is not associated with cancer outside of the hematopoietic stem cell compartment. A third set of germline loci are associated with the acquisition of CHIP mutations in specific driver genes. This previously was described in the *JAK2* 46/1 haplotype leading to *JAK2* p.V617F via a *cis* haplotype effect.<sup>26–28</sup> We now identify a novel *DNMT3A* CHIP specific locus at the *TCL1A* promoter specifically associated with increased risk of *DNMT3A* CHIP, but not other CHIP subsets.

We identify a convergence of common and rare germline genetic predisposition to leukocyte telomere length, MPN, large scale somatic chromosomal mosaicism and CHIP, suggesting shared causal mechanisms. Importantly, to date, only CHIP with leukemogenic driver mutations (as opposed to somatic chromosomal mosaicism<sup>9</sup> or CHIP with unknown driver mutations<sup>14</sup>) has been robustly associated with non-oncologic diseases independently of age. The partially overlapping genetic predisposition we observe across these three clonal phenomena suggest that although there may be similar genetic architecture that predispose individuals to acquiring a somatic mutation, the specific change may be particularly relevant to atherosclerotic disease as opposed to the general phenomenon of clonal hematopoiesis itself.



Third, our work underscores the benefits of studying genomes from individuals of diverse ancestries. The inclusion of a significant number of African Ancestry samples in TOPMed permitted the discovery of the *TET2* locus which was not present in other ancestries. Further inclusion of diverse individuals in genomic analyses is likely to highlight additional new biological pathways.

Important limitations of our study include reduced sensitivity for detecting CHIP with low allele fractions (VAF 2–5%) even with high-coverage whole genome sequencing. Ultrasensitive targeted sequencing can facilitate detection of such leukemogenic mutations at exceedingly low VAFs but the clinical consequences of this much more pervasive phenomenon, as well as determinants of progression to CHIP is not well understood currently.<sup>29</sup> Furthermore, the cross-sectional analyses of CHIP with non-genetic risk factors and biomarkers limit conclusions regarding temporal relationships between CHIP and these features; however, these observations still permit risk prediction for CHIP presence. Notably, inflammatory biomarker analyses are concordant with prior model experiments indicating elevations of observed inflammatory biomarkers as a consequence of CHIP.<sup>5,10</sup> Lastly, given the age-dependence of CHIP, it is likely that many individuals not observed to have CHIP in this study will develop CHIP in the future.

Overall, comprehensive simultaneous germline and somatic analyses of blood-derived whole genome sequence data demonstrates that germline variation influences the acquisition of somatic mutations in blood cells. Importantly, we anticipate that the TOPMed CHIP dataset defined here will be a valuable tool in establishing associations of CHIP with diverse heart, lung, blood and sleep traits.

## Methods

### Study Samples

Whole genome sequencing (WGS) was performed on 97,691 samples sequenced as part of 52 studies contributing to the NHLBI TOPMed research program Freeze 6 release as previously described for discovery analyses.<sup>7</sup> An additional distinct set of 9,389 WGS sequenced samples from the NHLBI TOPMed Freeze 8 release were used for replicating the *TET2* genetic association. Study designs include prospective cohorts, families, population isolates, and case-only collections. A subset of the studies focus on heart (~40%) or lung (~30%) phenotypes, with the remainder representing prospective population cohorts or electronic health record linked cohorts which have been assessed for many diverse phenotypes. None of the studies which comprise TOPMed selected individuals for sequencing on the basis of hematologic malignancy. Approximately 82% of participants are U.S. residents with diverse ancestry and ethnicity (40% European, 32% African, 16% Hispanic/Latino, 10% Asian). Each of the constituent studies provided informed consent on the participating samples. Details on participating cohorts and samples is provided in Supplementary Table 1. The age of participants at time of blood draw was obtained for a subset of 82,807 of the samples. The median age was 55, the mean age 52.5, and the maximum age 98. The age distribution varied across the constituent cohorts.

Written informed consent was obtained from all human participants by each of the studies that contributed to TOPMed with approval of study protocols by ethics committees at participating institutions as summarized in Supplementary Table 1. Each study received institutional certification prior to deposition in dbGaP which certified that all relevant institutional ethics committees approved the individual studies and that the genomic and phenotypic data submission was compliant with all relevant ethical regulations. This certification was deposited in dbGaP along with the data. Secondary analysis of the TOPMed dbGaP data as described in this manuscript was approved by the Partners Healthcare Institutional Review Board. All relevant ethics committees approved this study and this work is compliant with all relevant ethical regulations.

### WGS Processing, Variant Calling and CHIP annotation

BAM files were remapped and harmonized through a previously described unified protocol.<sup>30</sup> SNPs and short indels were jointly discovered and genotyped across the TOPMed samples using the GotCloud pipeline.<sup>31</sup> An SVM filter was trained to discriminate between true variants and low-quality sites. Sample quality was assessed through pedigree errors, contamination estimates, and concordance between self-reported sex and genotype inferred sex. Variants were annotated using snpEff 4.3.

Putative somatic SNPs and short indels were called with GATK Mutect2<sup>8</sup> (<https://software.broadinstitute.org/gatk>). Briefly, Mutect2 searches for sites where there is evidence for variation, and then performs local reassembly. It uses an external reference of recurrent sequencing artifacts termed a “panel of normal samples” to filter out these sites, and calls variants at sites where there is evidence for somatic variation. The panel of normal samples used for our study included 100 randomly selected individuals under the age of 40 years. Absence of a hotspot CHIP mutation was verified prior to inclusion in the panel of normal set. An external reference of germline variants<sup>32</sup> was provided to filter out likely germline calls. We deployed this variant calling process on Google Cloud using Cromwell (<https://github.com/broadinstitute/cromwell>). The caller was run individually for each sample with the same settings. The Cromwell WDL configuration file is available from the authors upon request.

Samples were annotated as having CHIP if the Mutect2 output contained one or more of a pre-specified list of putative CHIP variants as previously described<sup>2,5</sup> (Supplementary Table 2) at a VAF >2%.

### WGS sensitivity to detect CHIP

To empirically demonstrate the sensitivity of CHIP detection and VAF, we re-analyzed sequence data from 30 samples with CHIP from a previously published cohort.<sup>33</sup> These samples were sequenced to >400x depth. We bioinformatically down-sampled the reads to the range of sequencing depths compatible with whole exome and whole genome sequencing. The TOPMed WGS samples were sequenced to a median depth of ~40x, although sequencing of any particular region was typically 30x-50x. Across this range of sequencing depths we observe robust ability to call CHIP with VAF >10%, which is the most clinically actionable subset of CHIP. We also capture approximately half of the CHIP

calls in the VAF 5–10% range. To reliably capture CHIP in the 5–10% range requires ~100x sequencing depth commonly done in whole exome sequencing, but even at this sequencing depth the majority of the VAF 2–5% CHIP calls are not reliably detected. (Extended Data Figure 11)

### Amplicon sequencing validation

To evaluate the fidelity of our TOPMed WGS CHIP dataset, we performed technical validation of 76 CHIP mutations in 72 samples using targeted deep sequencing. All 76 of 76 CHIP mutations identified with WGS were also identified with targeted deep sequencing. CHIP mutations were validated by single-molecule molecular inversion probe sequencing (smMIPS).<sup>34</sup> Capture probes were designed to tile all coding exons (+/- 5 bp) for 12 of the mostly highly prevalent CHIP genes plus four recurrent mutation hotspots, totaling 44.5 kb. Probes were synthesized as a pool by CustomArray, Inc., amplified using Q5 DNA polymerase (NEB) using outer flanking primers, and digested with BbsI-HF (NEB) to remove adaptors. For each sample, captures were performed with 500 ng gDNA and converted to dual-barcoded Illumina sequencing libraries as described.<sup>35</sup> Sequence capture libraries were pooled for paired-end 150 bp sequencing on a HiSeq 4000 lane. Resulting reads were aligned with bwa mem and processed using the mimips pipeline (source code at <https://github.com/kitzmanlab/mimips>) to trim capture probe sequences, and to remove reads with duplicated unique molecular identifiers. Somatic variants were called by MuTect2 as described above and confirmed by manual inspection with IGV.

### Somatic Chromosomal Mosaic Detection

In order to assess the relationship between CHIP and clonal mosaicism reflecting chromosomal mutation, we sought to characterize large (megabase-scale) acquired chromosomal alterations leading to allelic imbalance using existing SNP array data on a subset of the samples in this analysis. To do so, we compared statistically reconstructed haplotypes (using MaCH<sup>36</sup>) with the patterns of “B allele” frequencies (BAFs), measured via SNP array. Regions of nonrandom similarities between the estimated haplotypes and BAFs were detected with hapLOH<sup>37</sup>, and indicate acquired chromosomal alterations. We identified genomic allelic imbalance events using a threshold of a posterior probability for allelic imbalance > 0.8 and event size > 1Mb. We excluded allelic imbalance events with fewer than ten markers and removed potential germline duplications if a detected event exhibited the following: 1) 50% reciprocal overlap with database of genomic variants (DGV) and 2) was not determined to be a deletion or LRR deviations > 0.08, size < 5Mb and BAF deviations > 0.1. Phasing and event detection was performed in SyQADA.<sup>38</sup>

### Blood traits

Conventionally measured 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. Phenotypes were collected by each cohort, centrally harmonized by the TOPMed Data Coordinating Center (DCC). Additional documentation about harmonization algorithms for each specific trait is available

from the TOPMed DCC and accompanies the data on the dbGaP TOPMed Exchange area. Up to 37,653 individuals from 10 cohorts were utilized for this analysis that had one or more blood traits measured concurrently or following the blood draw used for CHIP ascertainment. Traits were first log<sub>2</sub> normalized and then analyzed using a general linear regression model with CHIP status, age, sex, study and the first 10 ancestry principal components as covariates.

### Lipid phenotypes

Conventionally measured plasma lipids, including total cholesterol, LDL-C, HDL-C, and triglycerides, were included for analysis. LDL-C was either calculated by the Friedewald equation when triglycerides were <400 mg/dl or directly measured. Given the average effect of statins, when statins were present, total cholesterol was adjusted by dividing by 0.8 and LDL-C by dividing by 0.7. Triglycerides were natural log transformed for analysis. Phenotypes were harmonized by each cohort and deposited into dbGaP TOPMed Exchange area as previously described.<sup>39</sup> Up to 28,310 individuals from 19 cohorts were utilized for this analysis that had one or more lipid trait measured concurrently or following the blood draw used for CHIP ascertainment. Lipid traits were first normalized for age, sex and ancestry principal components and then analyzed using a general linear regression model with CHIP status, age, sex, study and the first 10 ancestry principal components as covariates.

### Inflammatory Markers

A set of markers previously implicated in mediating cardiometabolic disease were analyzed including: CD-40, CRP, E-Selectin, ICAM-1, IL-1b, IL-6, IL-10, IL-18, 8-epi PGF2a, Lp-PLA2 mass and activity, MCP1, MMP9, MPO, OPG, P-selectin, TNF-Alpha, TNF-Alpha Receptor 1, TNF-receptor 2. Phenotypes were collected by each cohort, centrally harmonized by the TOPMed DCC and then deposited into dbGaP TOPMed Exchange area. Additional documentation about harmonization algorithms for each specific trait is available from the TOPMed DCC and accompanies the data on dbGaP. Up to 22,092 individuals from 10 cohorts were utilized for this analysis that had one or more inflammatory marker measured concurrently or following the blood draw used for CHIP ascertainment. Inflammatory markers were first normalized using a log<sub>2</sub>(x+1) transformation and then analyzed using a general linear regression model with CHIP status, age, sex, study and the first 10 ancestry principal components as covariates.

### Mutational Signatures

We identified all putatively somatic singleton mutations in a subset of the TOPMed samples that included 3,764 cases with a single CHIP driver mutation and a randomly sampled set of 5,000 controls. Variants were filtered to ensure a depth  $\geq 25$  reads, a VAF < 35% and no overlap with the germline variant site list from TOPMed Freeze 5 (available: <https://bravo.sph.umich.edu/freeze5/hg38/>). Multiallelic variants and indels were excluded. We used [https://cancer.sanger.ac.uk/cosmic/signatures\\_v2](https://cancer.sanger.ac.uk/cosmic/signatures_v2) as a reference for mutation signatures and the MutationalPatterns R package to estimate the contributions of the signatures.<sup>13,40,41</sup> We defined a signature as being “differentially observed” if at least 99% of its observations

are in CHIP cases, or if at most 1% of its observations are in cases (i.e., one of cases or controls contains at least 99% of the signature observations).

### Single Variant Association

Single variant association for each variant in Freeze 8 with  $MAF > 0.1\%$  and  $MAC > 20$  was performed with SAIGE<sup>42</sup>, and analysis was performed using the TOPMed Encore analysis server (<https://encore.sph.umich.edu>). CHIP driver status was dichotomized into a case-control phenotype based on the presence of at least one driver mutation. Prior to running single variant association tests, a logistic mixed model was fit using the lme4 R package<sup>43</sup> to estimate the probability of the CHIP case control status conditional on a spline transformation of the centered age, genotype inferred sex, and cohort. The cohort was included as a random intercept which represents study specific contributions to the log-odds of CHIP at the mean sample age. Age was modeled with a spline to capture the non-linearity of the relationship between age and CHIP. This model was chosen over comparable models based on its AIC. Combining the age, inferred sex, and study into a single quantity aided the convergence of SAIGE compared to the inclusion of these terms separately. The first 10 principal components were also included as covariates.

Given that CHIP is unlikely to manifest in younger individuals, these individuals are effectively censored in our analysis set – that is, a young individual that does not presently have CHIP may still develop CHIP in the future. To avoid the power loss associated with misclassification of controls, we pruned these individuals from our analysis set. The single variant association analysis was run on a pruned set of samples that excluded those which had less than a 1% probability CHIP as estimated by the aforementioned model. This excluded 21,712 samples leading to a final analysis set of 65,405 which was used for downstream association analyses.

### Fine mapping

We applied FINEMAP 1.3<sup>44</sup> to the summary statistics from SAIGE, using the z-score and LD matrices as input. We fine-mapped the TET2 locus using the summary statistics from the African ancestry single variant summary statistics and estimated LD on the same set of samples using Plink 1.9. We set the maximum number of causal SNPs in the region to 10 and used a shotgun stochastic search.

### Transcriptome-wide association analysis

Multi-tissue gene expression and eQTL data were retrieved from the Genotype-Tissue Expression (GTEx) project (<https://www.gtexportal.org>). We applied the unified test for molecular signatures (UTMOST)<sup>20</sup> to perform cross-tissue transcriptome-wide association analysis for CHIP. We used cross-tissue gene expression imputation models trained from 44 tissues in GTEx. Gene-level association meta-analysis was performed using the generalized Berk-Jones test implemented in UTMOST (<https://github.com/Joker-Jerome/UTMOST>). Statistical significance was determined using a Bonferroni corrected p-value cutoff  $2.9 \times 10^{-6}$ .

## Rare Variant Analyses

Collapsing burden tests were applied to specific variant grouping schemes using EPACTS (<https://genome.sph.umich.edu/wiki/EPACTS>). The same covariates as the single variant tests were used on the same set of samples. We used burden tests due to their limited compute requirements, which were considerable for the number of variants and samples tested. Two grouping schemes were specified: the first groups coding variation, and the second groups putative regulatory elements in a relevant cell line. The first used all putative LOF variants as identified by snpEff. Given that some variants were present in both the Mutect2 calls and the germline variant calls, we pruned the LOF variants to exclude variants that were present in both call sets. The second grouping scheme used all variants in regions that were predicted enhancers for CD34 cells that had CADD scores of at least 10. Predicted enhancers were identified by the activity-by-contact model.<sup>45</sup>

## Predicting enhancer-gene regulation for TET2.

We used the Activity-by-Contact (ABC) model<sup>46</sup> to predict which enhancers regulate which genes in CD34+ hematopoietic progenitor cells, with minor modifications as follows.

Briefly, this model predicts the effect of each putative regulatory element (defined as a DNase peak within 5Mb of a given promoter) by multiplying the Activity of each element (estimated from DNase-seq and H3K27ac ChIP-seq) by its Contact with a target promoter (estimated from Hi-C data). The ABC score of a single element on a gene's expression is the predicted effect of that element divided by the sum of the predicted effects of all elements for a given gene.

We identified putative regulatory elements by using MACS2 to call peaks in DNase-seq data from mobilized CD34+ hematopoietic progenitor cells from the Roadmap Epigenome Project (downloaded from <http://egg2.wustl.edu/roadmap/data/byFileType/alignments/consolidated/E050-DNase.tagAlign.gz>) Initially we considered all peaks with p-value < 0.1. To further refine this list, we kept the 100,000 peaks with the highest number of DNase-seq reads. We then resized these peaks to be 500 bp in length centered on the peak summit, merging any overlapping peaks, and removed any peaks overlapping ENCODE "blacklisted regions"<sup>47</sup> (regions of the genome previously observed to accumulate anomalous numbers of reads in epigenetic sequencing experiments; downloaded from <https://sites.google.com/site/anshulkundaje/projects/blacklists>). To this peak list, we added 500 bp regions centered on the transcription start site of all genes. Any overlapping regions resulting from these additions or extensions were merged.

Within each putative regulatory element, we estimated enhancer Activity as the geometric mean of read counts from DNase-seq and H3K27ac ChIP-seq data from the Roadmap Epigenome Project (<http://egg2.wustl.edu/roadmap/data/byFileType/alignments/consolidated/E050-DNase.tagAlign.gz>, and [E050-H3K27ac.tagAlign.gz](http://egg2.wustl.edu/roadmap/data/byFileType/alignments/consolidated/E050-H3K27ac.tagAlign.gz)).

We estimated enhancer-promoter Contact from the KR-normalized Hi-C contact maps in primary CD34+ cells. We then calculated effect of each putative enhancer-gene connection by multiplying the Activity and Contact for that element and gene. Dividing the effect of each element by the sum of effects for all elements for a given gene yields the ABC score:



$$\text{ABC score}_{E-G} = \frac{A_E \times C_{E-G}}{\sum_{e \text{ within } 5 \text{ Mb}} A_e \times C_{e-G}}$$

To call predicted enhancer-gene connections, we used a threshold on the ABC score of 0.015. The rs79901204 variant overlapped an enhancer with ABC score of 0.0308 for TET2, which, based on comparison of ABC scores to large-scale enhancer perturbation datasets, corresponds to a positive predictive value of approximately 61%.

### Functional Evaluation of TET2 locus

The genomic region containing risk and non-risk allele of the variant rs79901204 (600bp) was synthesized as gblocks (IDT Technologies) and cloned into the Firefly luciferase reporter constructs (pGL4.24) using NheI and EcoRV sites. The Firefly constructs (500ng) were co-transfected with pRL-SV40 Renilla luciferase constructs (50ng) into 100,000 K562 cells (ATCC) using Lipofectamine LTX (Invitrogen) according to manufacturer's protocols. Cells were harvested after 48 hours and the luciferase activity measured by Dual-Glo Luciferase Assay system (Promega). K562 cell identity was validated using STR analysis. Mycoplasma testing was routinely performed on all cells used in the study and confirmed to test negative.

### CRISPR/Cas9 editing of CD34<sup>+</sup> human HSPCs

Editing of *TET2* enhancer and *TET2* coding regions was performed by electroporation of Cas9 Ribonucleoprotein complex (RNP) into CD34<sup>+</sup> human HSPCs. CD34<sup>+</sup> HSPCs from adult donors obtained from the Fred Hutchinson Cancer Research Center, Seattle, USA were thawed 24 hours prior to electroporation and cultured in HSC expansion conditions throughout the experiment (Stemspan II media with CC100 cytokine cocktail from Stem Cell Technologies and TPO (50ng/ul) and small molecule UM171 (35nM)). The RNP complex was made by mixing Cas9 (50 pmol) and modified sgRNAs from Synthego (100 pmol in total). HSPCs ( $3.75 \times 10^5$ ) resuspended in 20  $\mu$ l P3 solution were mixed with RNP and transferred to Nucleocuvette strips for electroporation with program DZ-100 (Lonza 4D Nucleofector). *TET2* gene expression was measured at 6 days post-electroporation.

For enhancer deletion experiments two guides targeting 5' and 3' ends of the enhancer element was used simultaneously (ENH\_sgRNA\_1: GGATTCTGTATTCGTCTGTG & ENH\_sgRNA\_2: TCTACTCACAGGGCCCAATG). For *TET2* coding disruption experiments single guides were used (TET2\_CDS1: TGGAGAAAGACGTAACCTCG & TET2\_CDS2: TCTGCCCTGAGGTATGCGAT). For negative control, a guide targeting *AAVS1* site was used (GGGGCCACTAGGGACAGGAT). Editing efficiency of *TET2* CDS and *AAVS1* guides were measured by Sanger sequencing followed by TIDE analysis. Editing efficiency of *TET2* enhancer deletion was measured by PCR and agarose gel electrophoresis.

### Colony-forming unit cell assays

3 days post RNP-electroporation, 500 CD34<sup>+</sup> HSPCs were plated in 1ml methylcellulose media (# H4034, Stem Cell Technologies). Primary CFU-C colonies were counted after 14



days. For the colony replating experiments, 2 weeks after the primary plating, the colonies from three plates were pooled, washed with PBS, and the cells were plated in new methylcellulose media at 25,000 cells/ml for an additional 2 weeks.

### RNA-Sequencing and eQTL Analysis:

RNA-Sequencing was performed on peripheral blood mononuclear cells from a subset of the MESA cohort. Alignment to the GRCh38 reference genome was done using STAR 2.5.3a.<sup>48</sup> Gene Quantification and quality control was performed using RNA-SeQC 1.19.<sup>49</sup> For RNA-SeQC, isoforms were collapsed into a single transcript per gene using the procedure described at [https://github.com/broadinstitute/gtex-pipeline/blob/master/gene\\_model/](https://github.com/broadinstitute/gtex-pipeline/blob/master/gene_model/). Samples that failed the RNA-Seq QC, fingerprinting, or expression-based sex check were filtered out. Further details on the RNASeq pipeline are available here: [https://www.nhlbiwgs.org/sites/default/files/TOPMed\\_RNAseq\\_pipeline\\_COREyr2.pdf](https://www.nhlbiwgs.org/sites/default/files/TOPMed_RNAseq_pipeline_COREyr2.pdf)

Analysis was performed using samples from 247 African Americans from MESA cohort Exam 1. Transcript expression was converted to TPM units (transcripts per million) and log2-transformed for analysis consistent with the GTEx consortium<sup>50</sup> best practices. Analysis of rs79901204 with *TET2* expression was performed using a linear mixed model adjusting for age at blood draw, sex, PC1–10 of population stratification from the WGS data, sequencing batch, and kinship relatedness matrix.

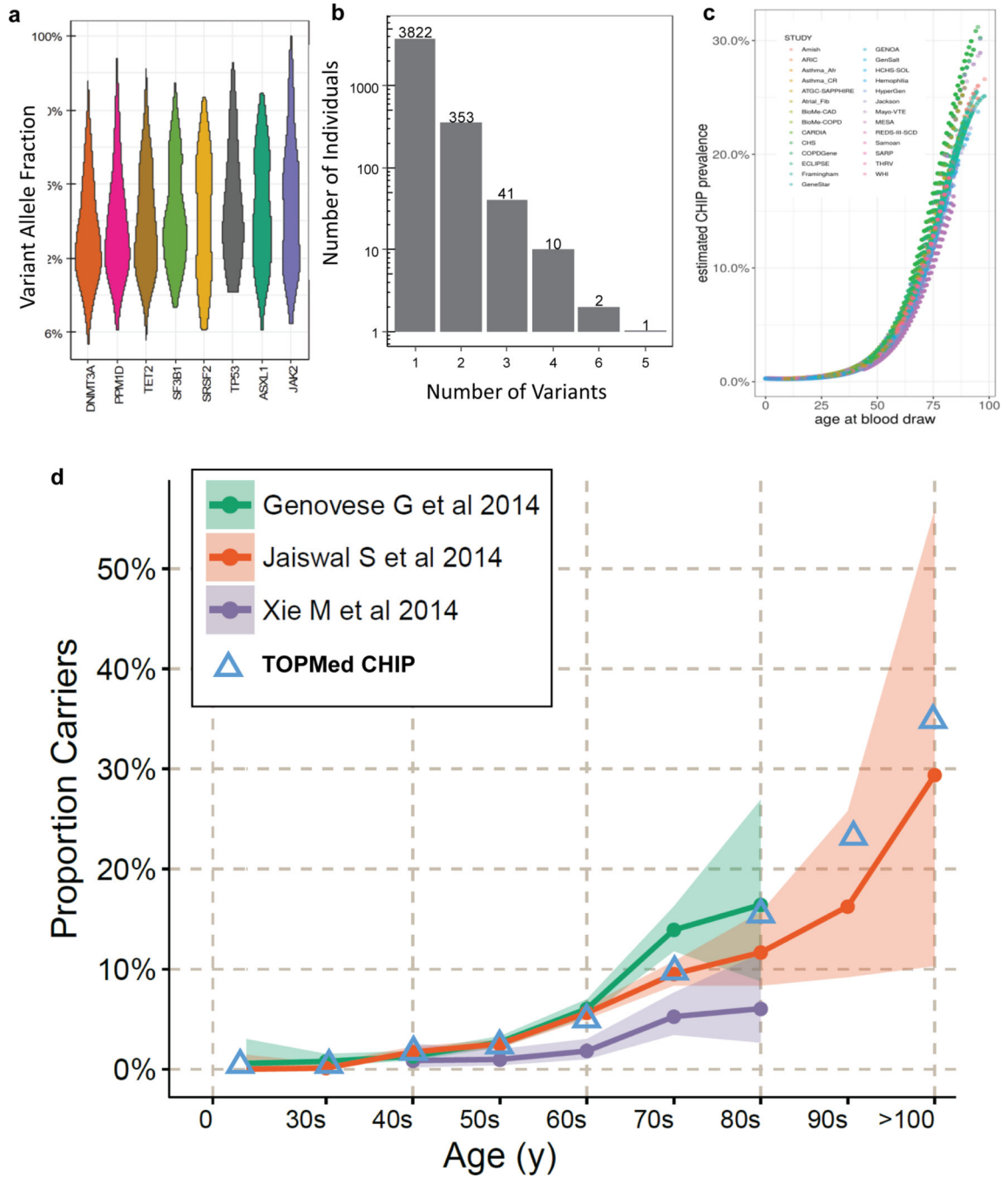
### Genome-wide Methylation-QTL analysis of *TET2* risk locus

Illumina Methylation EPIC 850K array data interrogating over 850,000 CpG DNA methylation sites was generated at the University of Washington's Northwest Genomic Center from blood samples collected from African Americans at the Jackson Heart Study baseline exam. Fluorescent signal intensities were preprocessed with the R package *minfi*<sup>51</sup> using the normal-exponential out-of-band (noob) background correction method with dye-bias normalization. N = 1747 total samples (1097 women and 650 men) remained after severe outliers were identified and removed. 71 individuals were positive for CHIP and 100 were carriers of the rs79901204 variant.

Methylation levels at each CpG site were then quantified as  $\beta$  values, defined as the ratio of intensities between methylated (M) and unmethylated (U) signals where  $\beta = M/(M+U+100)$ . Values therefore ranged from  $\beta = 0$  (completely unmethylated) to  $\beta = 1$  (completely methylated). Batch correction for assay plate position was performed on the  $\beta$  values via *ComBat*.<sup>52</sup> Relative leukocyte cell counts (CD8+ T-lymphocytes, CD4+ T-lymphocytes, Natural Killer cells, B cells, Monocytes, and Granulocytes) were estimated as previously described by Houseman<sup>53</sup> and Horvath<sup>52</sup>.

To investigate methylation in the *TET2* locus, a linear mixed effects model was fitted using CpGassoc<sup>53</sup> in R 3.6.0 with rs79901204 as the predictor and the batch-corrected methylation  $\beta$  levels as the dependent variable, adjusting for age, sex, estimated cell counts, the top 10 principal components of genetic ancestry, and CHIP status. A Bonferroni corrected threshold of  $p = 5.8 \times 10^{-8}$  was used to establish statistical significance.

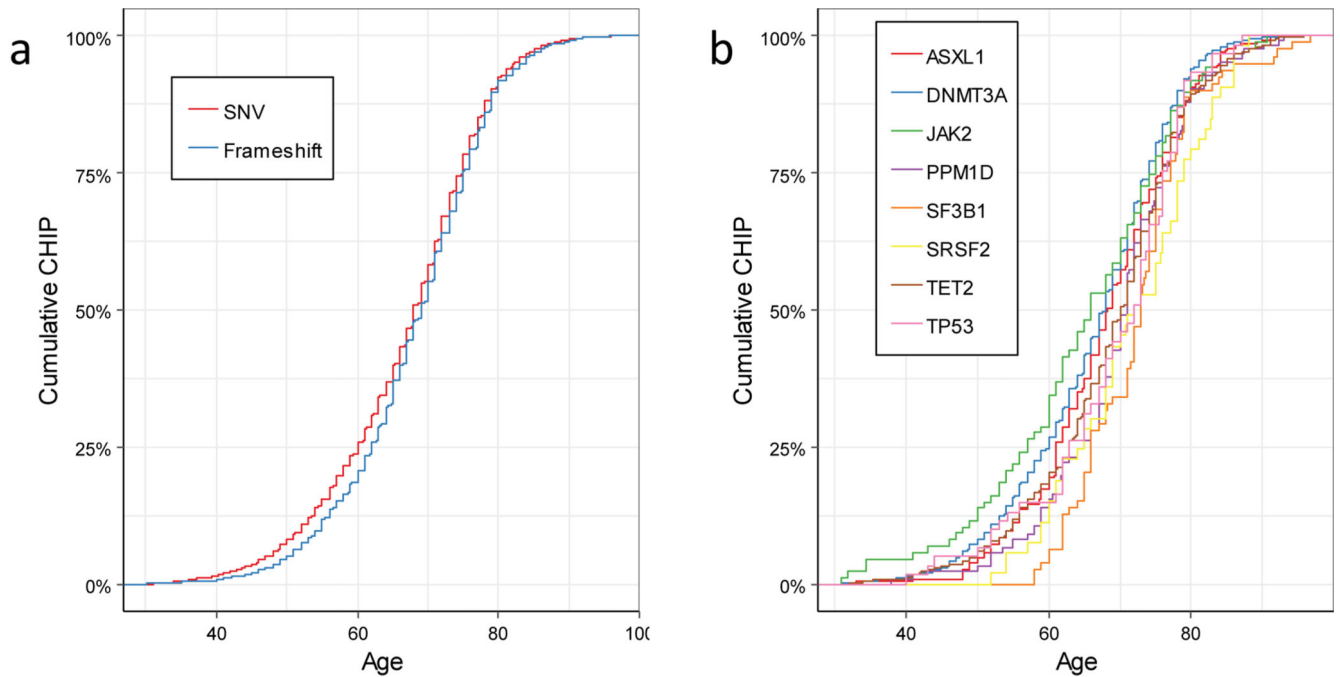
Extended Data



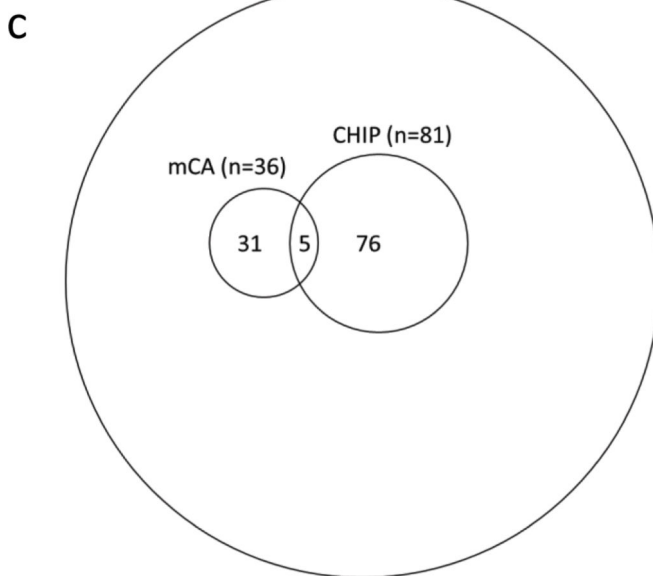
Extended Data Fig. 1|. Characterizing TOPMed CHIP.

a, There was marked heterogeneity of CHIP clone size as measured by variant allele fraction by CHIP driver gene. Violin plot spanning minimum and maximum values calculated on full dataset (Supplementary Table 3). Sample size for each element in violin plot displayed in Fig. 1, b, 90% of individuals with CHIP had only one CHIP driver mutation identified c, CHIP prevalence with age was highly concordant across sequenced cohorts. CHIP

prevalence was estimated from a logistic mixed model with spline-transformed age, sex, and cohort included as predictors. The cohort was included as a random intercept. Sample size for each cohort listed in Supplementary Table 1. **d**, CHIP prevalence with age in this study (blue triangles,  $N=82,807$ ) was highly consistent with previously observed CHIP prevalence (dots represent mean point prevalence with shaded area represents 95% confidence interval;  $N_{\text{Genovese}}=12,380$ ;  $N_{\text{Jaiswal}}=17,182$ ;  $N_{\text{Xie}}=2,728$ ).

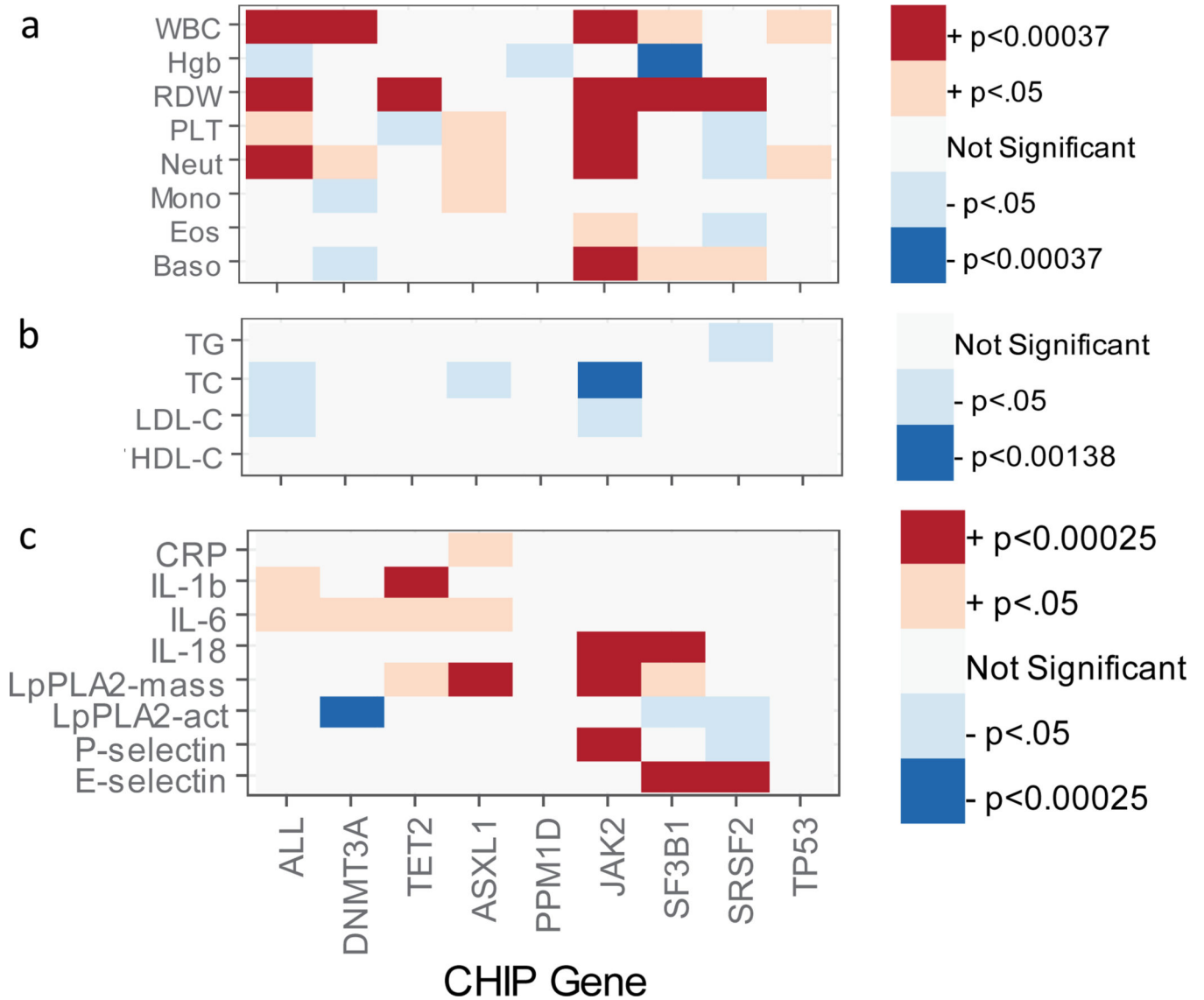


WHI: Whole genome sequencing & array genotyping (n=855)



**Extended Data Fig. 2|. CHIP age association by mutational mechanism, gene and overlap with somatic chromosomal mosaicism.**

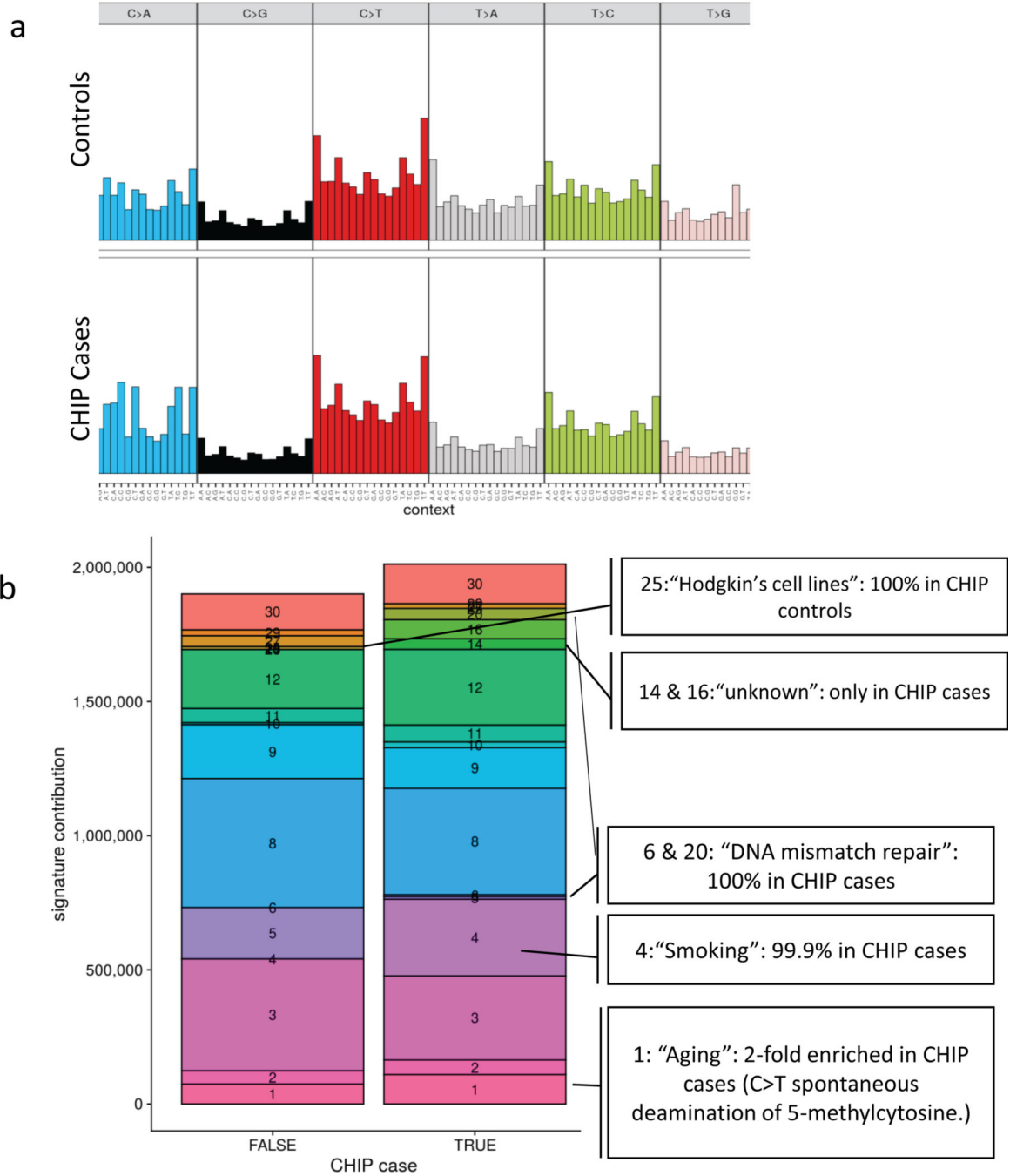
**a**, cumulative density plot of CHIP incidence with age stratified by single nucleotide variant (SNV) vs frameshift mutations. SNVs were observed in younger individuals than Frameshift mutations (N=4,939; two-sided wilcox rank sum test  $p=0.01$ ). **b**, cumulative density plot of CHIP incidence with age stratified by driver gene. **c**, 855 elderly WHI individuals (mean age: 70) with both whole genome and the array genotyping data available were interrogated for large-scale mosaic chromosomal rearrangements. The two somatic events did not occur more than would be expected by chance (hypergeometric  $p=0.25$ ).



**Extended Data Fig. 3|. CHIP associates with Blood, Lipid, and Inflammatory traits.**

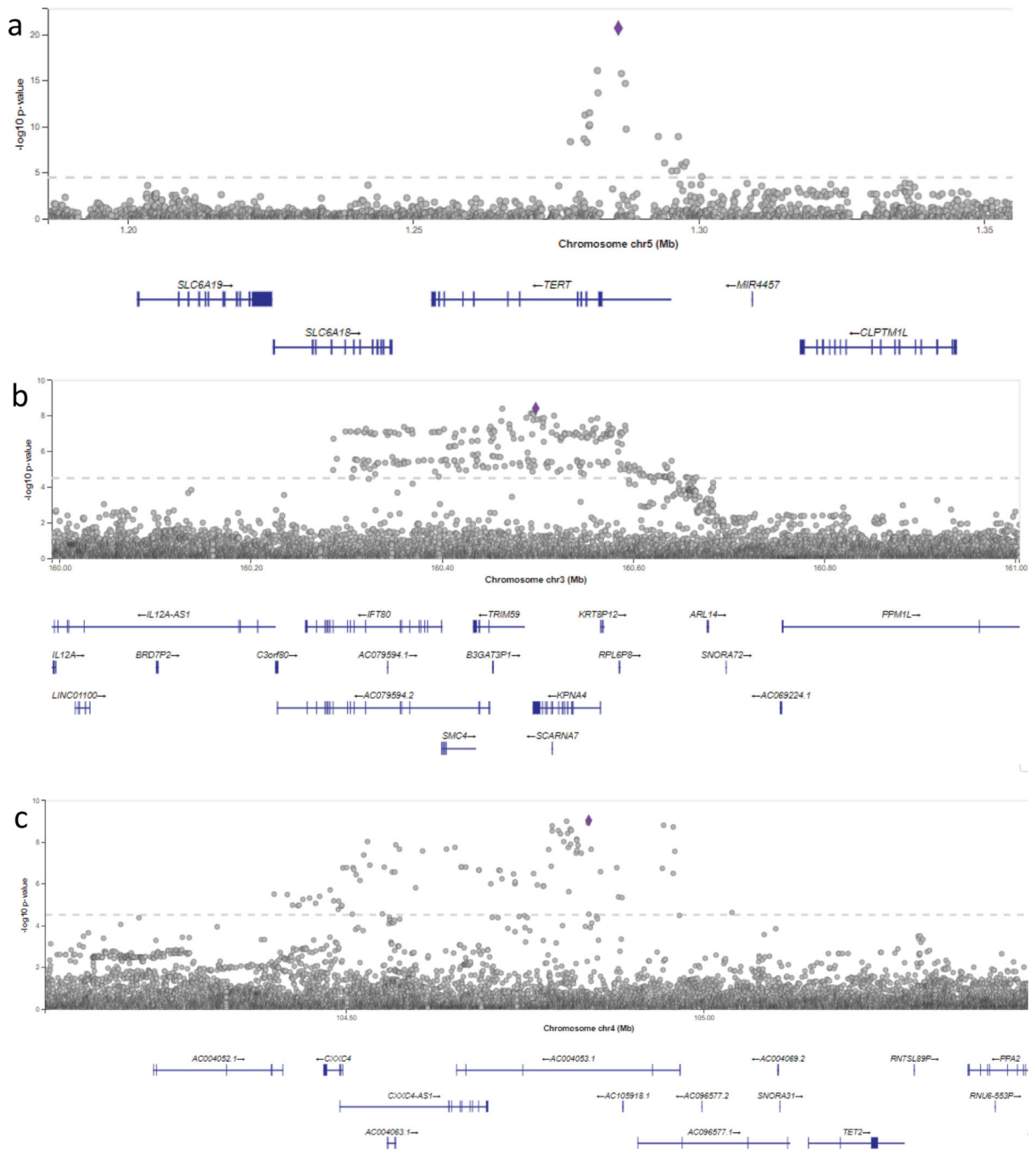
**a**, CHIP consistently associated with increased Red Cell Distribution Width (RDW). *JAK2*, *SF3B1* and *SRSF2* showed driver gene specific effects on blood traits (see Supplementary Table S5) **b**, CHIP status was not consistently associated with lipid traits, other than *JAK2* CHIP which was associated with decreased total cholesterol and a trend towards decreased LDL (see Supplementary Table S6) **c**, CHIP status is associated with inflammatory markers,

however notable heterogeneity existed across CHIP mutations (see Supplementary Table S7). Associations utilized a two-sided t-test from a multivariate general linear model including age, smoking, race and gender and study center and were not adjusted for multiple comparisons. Sample sizes and exact p-values for each phenotype are listed in Supplementary Tables 5–7.



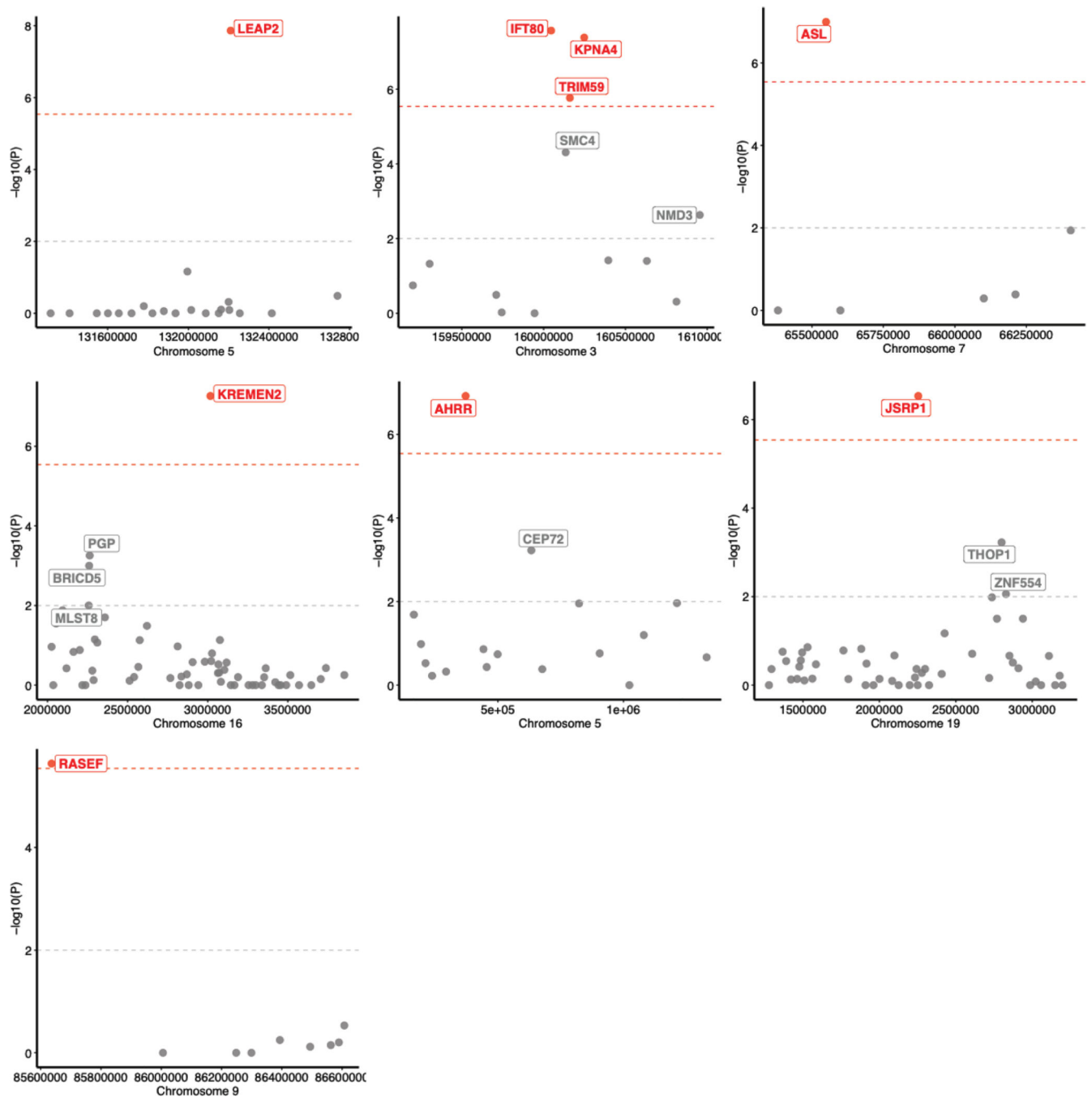
Extended Data Fig. 4|. CHIP passenger somatic mutation spectrum.

**a**, Singleton mutation counts by nucleotide context in CHIP Cases and Controls. **b**, Signature contribution in CHIP cases and controls identified differential enrichment



**Extended Data Fig. 5]. CHIP Single variant association regional association plots.**

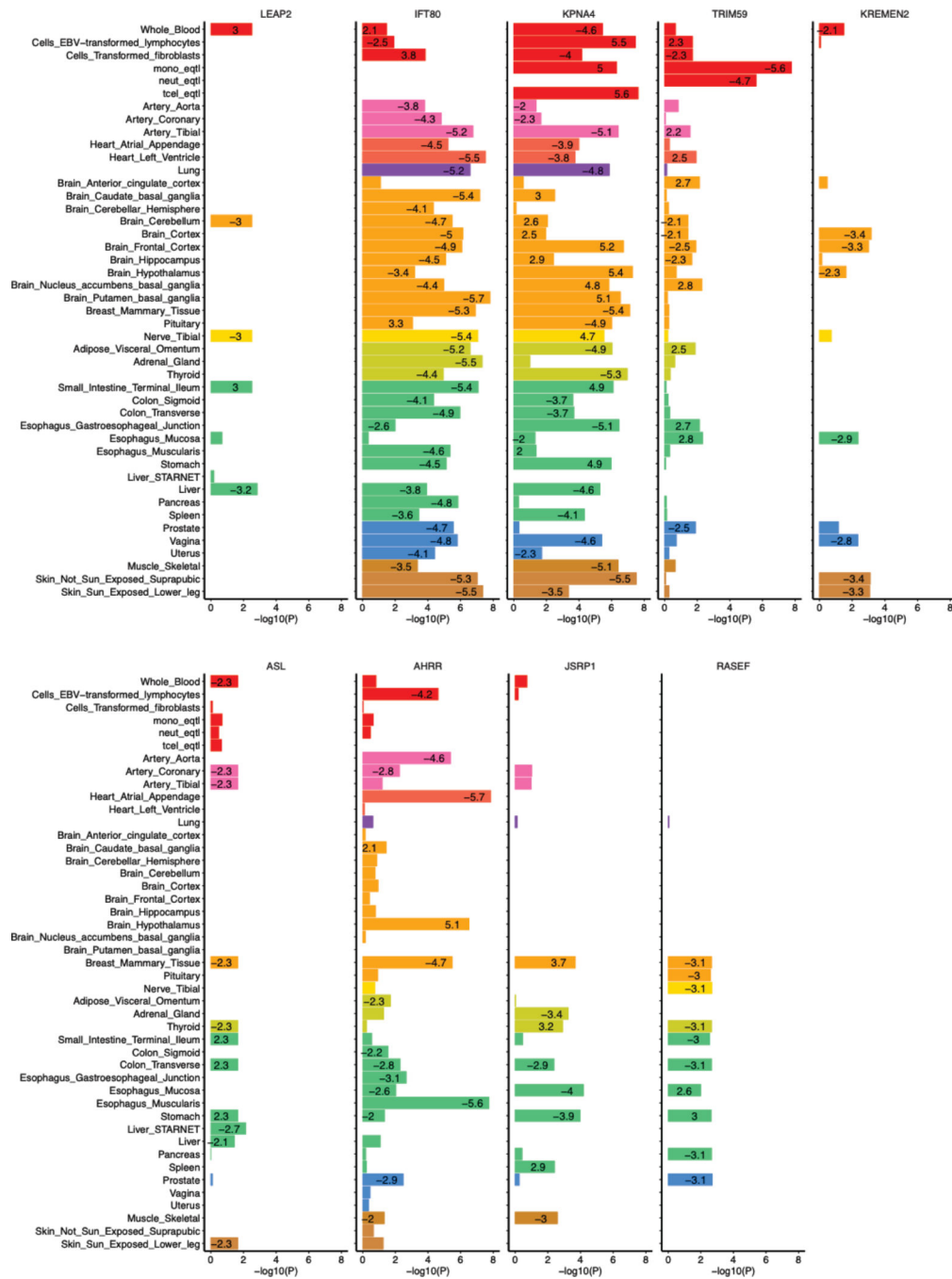
**a**, TERT locus **b**, TRIM59/KPNA4 locus **c**, TET2 locus. Two-sided association testing performed using SAIGE (N=65,405 individuals, see methods)



**Extended Data Fig. 6]. CHIP transcriptome-wide association study (TWAS) results across 48 tissues identified 7 significant loci.**

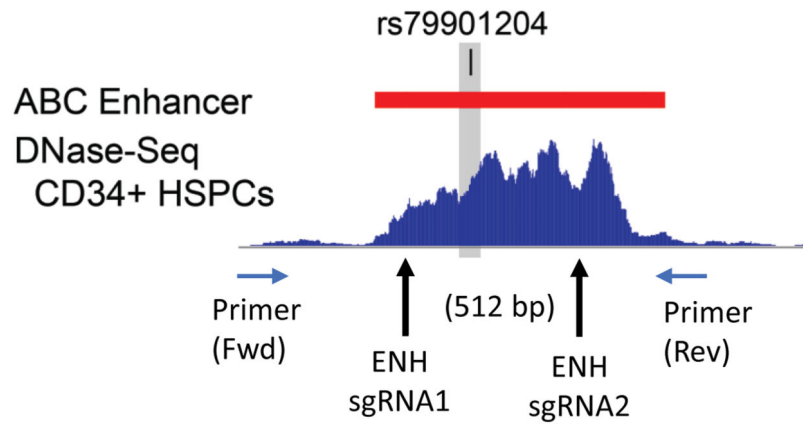
UTMOST algorithm applied to CHIP genome wide association study results from  $n=65,405$  individuals (see methods). Genomic coordinates listed on x-axis. P-value from generalized Berk-Jones test on Y axis. Multiple hypothesis corrected threshold,  $p < 2.9 \times 10^{-6}$  displayed as dotted red line.



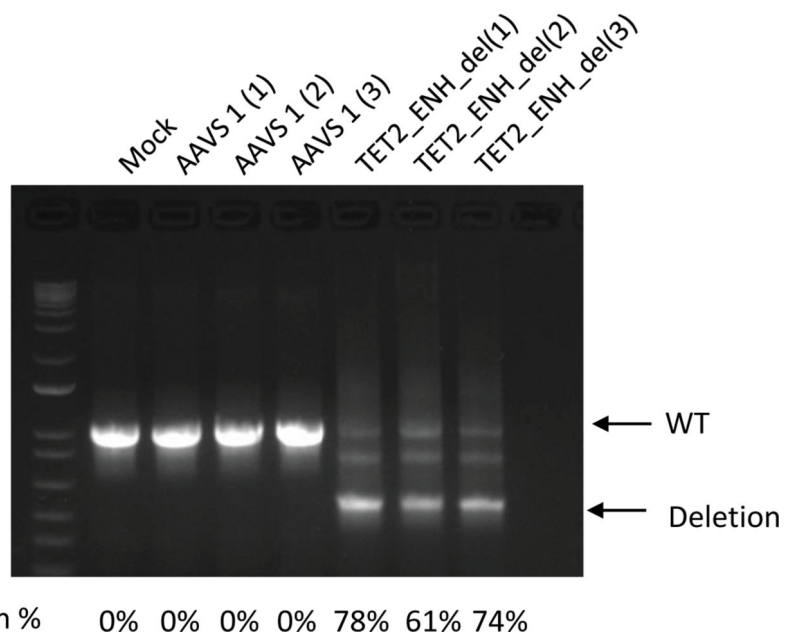


**Extended Data Fig. 7]. Tissue-specific results from the top 9 overall UTMOST-significant genes.** UTMOST algorithm applied to CHIP genome wide association study results from n=65,405 individuals. P-value from generalized Berk-Jones test. eQTL z-scores for associations with P<0.05 are displayed in each bar. GTEx eQTL tissue listed on Y-axis.

a

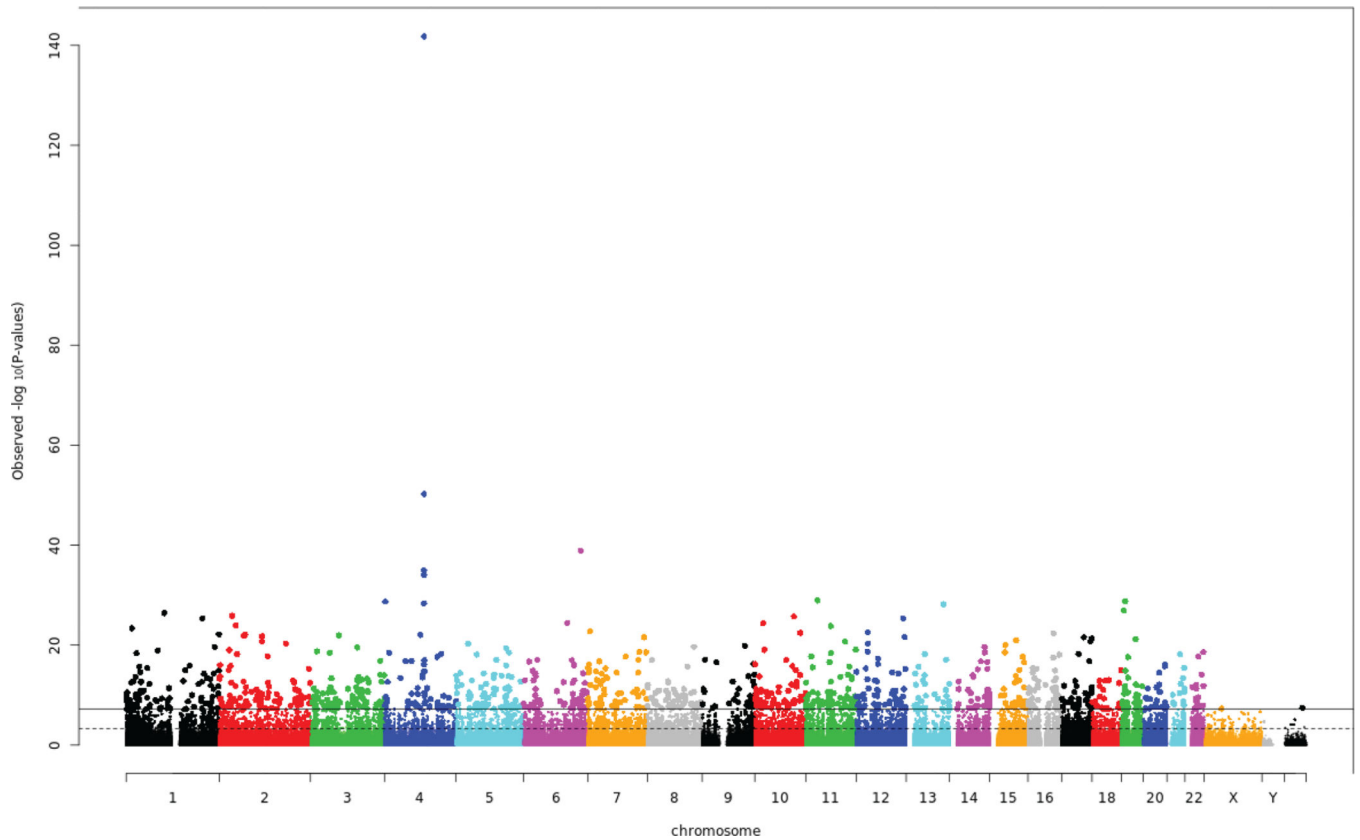


b

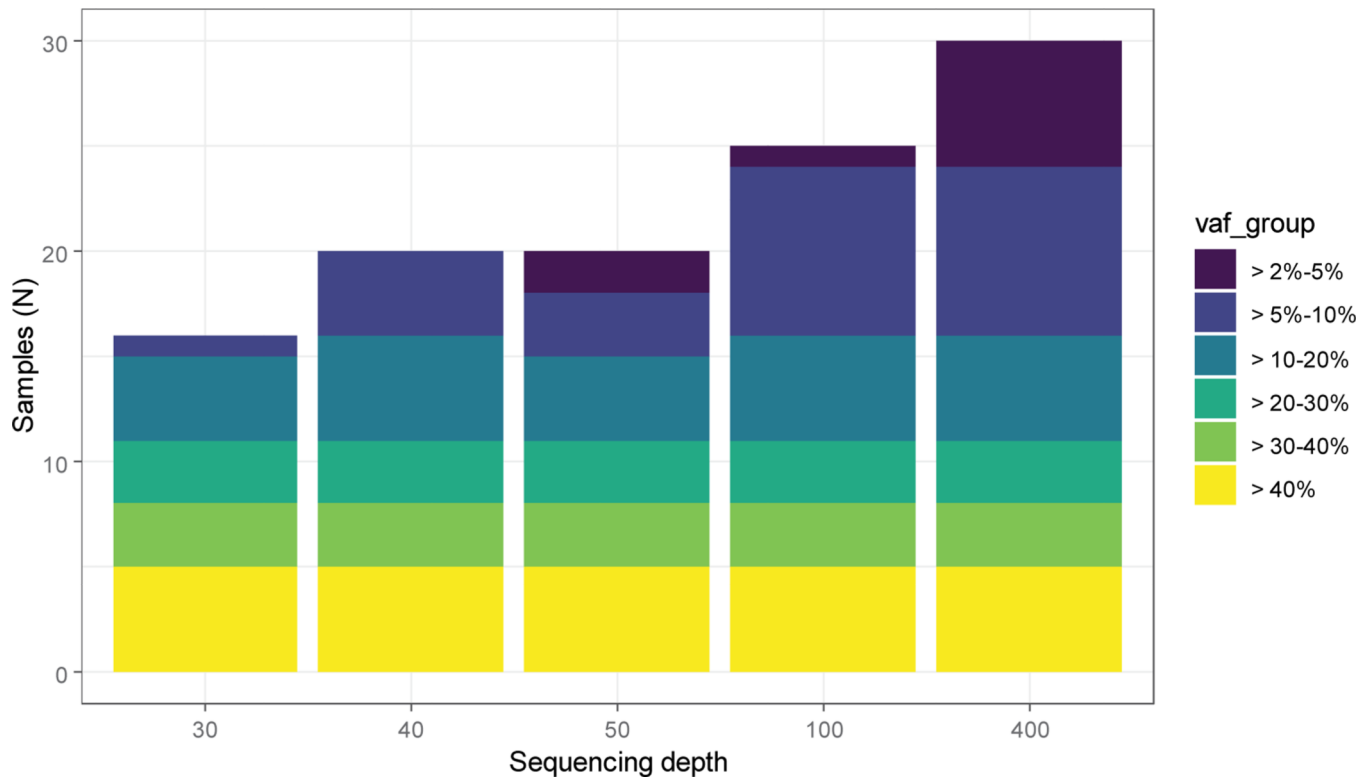


**Extended Data Fig. 8 |. CRISPR/Cas9 editing efficiency of TET2 Enhancer deletion in primary CD34+ HSPCs.**

a, Schematic showing the position of the two sgRNAs used to delete the TET2 enhancer (512bp) containing rs79901204. B, Gel electrophoresis image of PCR products from genomic DNA of edited HSPCs indicating unedited (WT) and deletion bands at sgRNA target site. Percentages of deletion alleles determined by band intensity and is shown below each lane. The experiment contains 3 biological replicates and was performed once.



**Extended Data Fig. 9 | rs79901204 associated with genome wide differential methylation signal,** Methylation Quantitative Trait association results of rs79901204 variant with cpg methylation probes identify an altered peripheral leukocyte methylation profile genome wide in  $N = 1747$  individuals. The strongest signal is at the chr4 TET2 locus. P-values on Y-axis derived from two-sided linear mixed effects model (see methods). To account for multiple hypothesis testing, a Bonferroni threshold of  $p < 5.8 \times 10^{-8}$  was used to establish statistical significance.



**Extended Data Fig. 10 |. Sensitivity of CHIP detection at various VAFs across sequencing depths.** A set of 30 samples from a previously published CHIP cohort (Gibbons et al, 2017) were computationally down sampled to 30x, 40x, 50x, 100x and 400x sequencing depth. TOPMed WGS data was typically in the 40x depth range across CHIP genes. WGS data has excellent sensitivity to detect CHIP clones with VAF >10%, and ~50% sensitivity to detect CHIP VAF 5–10%, with minimal ability to detect CHIP clones <5%.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Authors

Alexander G. Bick<sup>1,2,3,\*</sup>, Joshua S. Weinstock<sup>4,\*</sup>, Satish K. Nandakumar<sup>2,5</sup>, Charles P. Fulco<sup>2,6</sup>, Erik L. Bao<sup>2,5,7</sup>, Seyedeh M. Zekavat<sup>2,8</sup>, Mindy D. Szeto<sup>9,10</sup>, Xiaotian Liao<sup>2,5</sup>, Matthew J. Leventhal<sup>2</sup>, Joseph Nasser<sup>2</sup>, Kyle Chang<sup>11</sup>, Cecelia Laurie<sup>12</sup>, Bala Bharathi Burugula<sup>13</sup>, Christopher J. Gibson<sup>14</sup>, Amy E. Lin<sup>15</sup>, Margaret A. Taub<sup>16</sup>, Francois Aguet<sup>2</sup>, Kristin Ardlie<sup>2</sup>, Braxton D. Mitchell<sup>17,18</sup>, Kathleen C. Barnes<sup>9,19</sup>, Arden Moscati<sup>20</sup>, Myriam Fornage<sup>21,22</sup>, Susan Redline<sup>3,23,24</sup>, Bruce M. Psaty<sup>25,26,27,28</sup>, Edwin K. Silverman<sup>3,29</sup>, Scott T. Weiss<sup>3,29</sup>, Nicholette D. Palmer<sup>30</sup>, Ramachandran S. Vasan<sup>31</sup>, Esteban G. Burchard<sup>31,32</sup>, Sharon L. R. Kardia<sup>33</sup>, Jiang He<sup>34,35</sup>, Robert C. Kaplan<sup>36,37</sup>, Nicholas L. Smith<sup>26,28,38</sup>, Donna K. Arnett<sup>39</sup>, David A. Schwartz<sup>40</sup>, Adolfo Correa<sup>41</sup>, Mariza de Andrade<sup>42</sup>, Xiuqing Guo<sup>43</sup>, Barbara A. Konkle<sup>44,45</sup>, Brian Custer<sup>46,47</sup>, Juan M.

Peralta<sup>48</sup>, Hongsheng Gui<sup>49</sup>, Deborah A. Meyers<sup>50</sup>, Stephen T. McGarvey<sup>51</sup>, Ida Yii-Der. Chen<sup>52</sup>, M. Benjamin Shoemaker<sup>53</sup>, Patricia A. Peyser<sup>33</sup>, Jai G. Broome<sup>12</sup>, Stephanie M. Gogarten<sup>12</sup>, Fei Fei Wang<sup>12</sup>, Quenna Wong<sup>12</sup>, May E. Montasser<sup>17</sup>, Michelle Daya<sup>9</sup>, Eimear E. Kenny<sup>54</sup>, Kari E. North<sup>55</sup>, Lenore J. Launer<sup>56</sup>, Brian E. Cade<sup>23,57</sup>, Joshua C. Bis<sup>25</sup>, Michael H. Cho<sup>29</sup>, Jessica Lasky-Su<sup>3,29</sup>, Donald W. Bowden<sup>30</sup>, L. Adrienne. Cupples<sup>58</sup>, Angel C. Mak<sup>32</sup>, Lewis C. Becker<sup>59</sup>, Jennifer A. Smith<sup>33,60</sup>, Tanika N. Kelly<sup>34,35</sup>, Stella Aslibekyan<sup>61</sup>, Susan R. Heckbert<sup>26,28</sup>, Hemant K. Tiwari<sup>62</sup>, Ivana V. Yang<sup>40</sup>, John A. Heit<sup>63</sup>, Steven Lubitz<sup>2,3,64</sup>, Jill M. Johnsen<sup>44,45</sup>, Joanne E. Curran<sup>48</sup>, Sally E. Wenzel<sup>65</sup>, Daniel E. Weeks<sup>66</sup>, Dabeeru C. Rao<sup>67</sup>, Dawood Darbar<sup>68</sup>, Jee-Young Moon<sup>36</sup>, Russell P. Tracy<sup>69</sup>, Erin J. Buth<sup>12</sup>, Nicholas Rafaels<sup>19</sup>, Ruth J.F. Loos<sup>20,70</sup>, Peter Durda<sup>69</sup>, Yongmei Liu<sup>71</sup>, Lifang Hou<sup>72</sup>, Jiwon Lee<sup>23</sup>, Priyadarshini Kachroo<sup>3,29</sup>, Barry I. Freedman<sup>73</sup>, Daniel Levy<sup>74,75</sup>, Lawrence F. Bielak<sup>33</sup>, James E. Hixson<sup>76</sup>, James S. Floyd<sup>77</sup>, Eric A. Whitsetl<sup>78,79</sup>, Patrick T. Ellinor<sup>2,3,64</sup>, Marguerite R. Irvin<sup>61</sup>, Tasha E. Fingerlin<sup>80</sup>, Laura M. Raffield<sup>81</sup>, Sebastian M. Armasu<sup>82</sup>, Marsha M. Wheeler<sup>83</sup>, Ester C. Sabino<sup>84</sup>, John Blangero<sup>48</sup>, L. Keoki Williams<sup>49</sup>, Bruce D. Levy<sup>3,85</sup>, Wayne Huey-Herng Sheu<sup>86</sup>, Dan M. Roden<sup>87</sup>, Eric Boerwinkle<sup>88</sup>, JoAnn E. Manson<sup>3,89,90</sup>, Rasika A. Mathias<sup>58</sup>, Pinkal Desai<sup>91</sup>, Kent D. Taylor<sup>92</sup>, Andrew D. Johnson<sup>74,75</sup>, NHLBI Trans-Omics for Precision Medicine Consortium<sup>†</sup>, Paul L. Auer<sup>93</sup>, Charles Kooperberg<sup>94</sup>, Cathy C. Laurie<sup>12</sup>, Thomas W. Blackwell<sup>4</sup>, Albert V. Smith<sup>4</sup>, Hongyu Zhao<sup>95,96</sup>, Ethan Lange<sup>9</sup>, Leslie Lange<sup>9</sup>, Stephen S. Rich<sup>97</sup>, Jerome I. Rotter<sup>92</sup>, James G. Wilson<sup>98,99</sup>, Paul Scheet<sup>11</sup>, Jacob O. Kitzman<sup>13,100</sup>, Eric S. Lander<sup>2,101,102</sup>, Jesse M. Engreitz<sup>2,103</sup>, Benjamin L. Ebert<sup>2,3,14,104</sup>, Alexander P. Reiner<sup>26,94</sup>, Siddhartha Jaiswal<sup>105</sup>, Gonçalo Abecasis<sup>4,106</sup>, Vijay G. Sankaran<sup>2,3,5</sup>, Sekar Kathiresan<sup>2,3,107,108,#</sup>, Pradeep Natarajan<sup>2,3,64,#</sup>

## Affiliations

<sup>1</sup>Department of Medicine, Massachusetts General Hospital, Boston, MA.

<sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, MA.

<sup>3</sup>Harvard Medical School, Boston MA.

<sup>4</sup>Center for Statistical Genetics, Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, MI.

<sup>5</sup>Division of Hematology/Oncology, Boston Children's Hospital and Department of Pediatric Oncology, Dana-Farber Cancer Institute, Boston MA.

<sup>6</sup>Department of Systems Biology, Harvard Medical School, Boston, MA.

<sup>7</sup>Health Sciences and Technology Program, Harvard Medical School, Boston, MA.

<sup>8</sup>Yale School of Medicine, New Haven, CT.

<sup>9</sup>Division of Biomedical Informatics and Personalized Medicine, Department of Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO.

<sup>10</sup>Medical Scientist Training Program, University of Colorado Anschutz Medical Campus, Aurora, CO.

- <sup>11</sup>Department of Epidemiology, The University of Texas MD Anderson Cancer Center, Houston, TX.
- <sup>12</sup>Department of Biostatistics, University of Washington, Seattle, WA.
- <sup>13</sup>Department of Human Genetics, University of Michigan, Ann Arbor, MI.
- <sup>14</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA.
- <sup>15</sup>Division of Cardiovascular Medicine, Department of Medicine, Brigham and Women's Hospital, Boston, MA.
- <sup>16</sup>Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD.
- <sup>17</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD.
- <sup>18</sup>Geriatrics Research and Education Clinical Center, Baltimore Veterans Administration Medical Center, Baltimore, MD.
- <sup>19</sup>Colorado Center for Personalized Medicine, School of Medicine, University of Colorado, Aurora, CO.
- <sup>20</sup>Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY.
- <sup>21</sup>Brown Foundation Institute of Molecular Medicine, McGovern Medical School, University of Texas Health Science Center at Houston, Houston, TX.
- <sup>22</sup>Human Genetics Center, School of Public Health, University of Texas health Science Center at Houston, Houston, TX.
- <sup>23</sup>Division of Sleep and Circadian Disorders, Department of Medicine, Brigham and Women's Hospital, Boston, MA.
- <sup>24</sup>Department of Medicine, Beth Israel Deaconess Medical Center, Boston, MA.
- <sup>25</sup>Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA.
- <sup>26</sup>Department of Epidemiology, University of Washington, Seattle, WA.
- <sup>27</sup>Department of Health Services, University of Washington, Seattle, WA.
- <sup>28</sup>Kaiser Permanente Washington Health Research Institute, Seattle, WA.
- <sup>29</sup>Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital, Boston, MA.
- <sup>30</sup>Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC.
- <sup>31</sup>Departments of Medicine and Epidemiology, Boston University School of Medicine, Boston, MA.
- <sup>32</sup>Department of Medicine, University of California, San Francisco, CA.

<sup>33</sup>Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI.

<sup>34</sup>Department of Epidemiology, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA.

<sup>35</sup>Tulane University Translational Science Institute, New Orleans, LA.

<sup>36</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY.

<sup>37</sup>Fred Hutchinson Cancer Research Center, Division of Public Health Sciences, Seattle WA.

<sup>38</sup>Seattle Epidemiologic Information and Research Center, Department of Veterans Affairs, Office of Research & Development, Seattle WA.

<sup>39</sup>College of Public Health, University of Kentucky, Lexington KY.

<sup>40</sup>Department of Medicine, University of Colorado, Aurora, CO.

<sup>41</sup>Departments of Medicine and Population Health Science, University of Mississippi Medical Center, Jackson, MS.

<sup>42</sup>Department of Health Sciences Research, Mayo Clinic, Rochester, MN.

<sup>43</sup>Institute for Translational Genomics and Population Sciences, Department of Pediatrics, Los Angeles Biomedical Research Institute at Harbor-UCLA Medical Center, Torrance, CA.

<sup>44</sup>Bloodworks Northwest, Seattle, WA.

<sup>45</sup>Department of Medicine, University of Washington, Seattle, WA.

<sup>46</sup>Vitalant Research Institute, San Francisco, CA.

<sup>47</sup>University of California at San Francisco, San Francisco, CA.

<sup>48</sup>Department of Human Genetics and South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley School of Medicine, Brownsville, TX.

<sup>49</sup>Center for Individualized and Genomic Medicine Research, Department of Internal Medicine, Henry Ford Health System, Detroit, MI.

<sup>50</sup>Division of Genetics, Genomics and Precision Medicine, University of Arizona, Tucson, AZ.

<sup>51</sup>Department of Epidemiology and International Health Institute, Brown University School of Public Health, Providence.

<sup>52</sup>Medical Genetics, Los Angeles Biomedical Research Institute at Harbor-UCLA Medical Center, Los Angeles, CA.

<sup>53</sup>Division of Cardiology, Vanderbilt University Medical Center, Nashville, TN.

<sup>54</sup>Institute for Genomic Health, Icahn School of Medicine at Mount Sinai, New York, NY.



- <sup>55</sup>Department of Epidemiology, University of North Carolina, Chapel Hill, NC.
- <sup>56</sup>Laboratory of Epidemiology, Demography, and Biometry, Intramural Research Program, National Institute on Aging, Bethesda, MD.
- <sup>57</sup>Division of Sleep Medicine, Harvard Medical School, Boston, MA.
- <sup>58</sup>Departments of Biostatistics and Epidemiology, Boston University School of Medicine, Boston, MA.
- <sup>59</sup>GeneSTAR Research Program, Department of Medicine, Johns Hopkins School of Medicine, Baltimore MD.
- <sup>60</sup>Survey Research Center, Institute for Social Research, University of Michigan, Ann Arbor, MI.
- <sup>61</sup>Department of Epidemiology, University of Alabama at Birmingham, Birmingham, AL.
- <sup>62</sup>Department of Biostatistics, University of Alabama at Birmingham, Birmingham, AL.
- <sup>63</sup>Division of Cardiovascular Diseases, Mayo Clinic, Rochester, MN.
- <sup>64</sup>Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA.
- <sup>65</sup>Department of Environmental and Occupational Health, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA.
- <sup>66</sup>Departments of Human Genetics and Biostatistics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA.
- <sup>67</sup>Division of Biostatistics, Washington University School of Medicine, St. Louis, MO.
- <sup>68</sup>Division of Cardiology, University of Illinois at Chicago, Chicago, IL.
- <sup>69</sup>Department of Pathology and Laboratory Medicine, Larner College of Medicine, University of Vermont, Burlington, VT.
- <sup>70</sup>Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New York, NY.
- <sup>71</sup>Division of Cardiology, Department of Medicine, Duke University Medical Center, Durham, NC.
- <sup>72</sup>Department of Preventive Medicine, Northwestern Feinberg School of Medicine, Northwestern University, Chicago, IL.
- <sup>73</sup>Department of Internal Medicine, Section on Nephrology, Wake Forest School of Medicine, Winston-Salem, NC.
- <sup>74</sup>Framingham Heart Study, Framingham, MA.
- <sup>75</sup>Population Sciences Branch, National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, MD.

- <sup>76</sup>Department of Epidemiology, Human Genetics and Environmental Sciences, University of Texas Health Science Center at Houston School of Public Health, Houston, TX.
- <sup>77</sup>Departments of Medicine and Epidemiology, University of Washington, Seattle, WA.
- <sup>78</sup>Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina, Chapel Hill, NC.
- <sup>79</sup>Department of Medicine, School of Medicine, University of North Carolina, Chapel Hill, NC.
- <sup>80</sup>Center for Genes Environment and Health, National Jewish Health, Denver, CO.
- <sup>81</sup>Department of Genetics, University of North Carolina, Chapel Hill, NC.
- <sup>82</sup>Department of Health Sciences Research, Mayo Clinic, Rochester, MN.
- <sup>83</sup>Department of Genome Science, University of Washington, Seattle, WA.
- <sup>84</sup>Instituto de Medicina Tropical, Faculdade de Medicina da Universidade de Sao Paulo, Sao Paulo, Brazil.
- <sup>85</sup>Division of Pulmonary and Critical Care Medicine, Department of Medicine, Brigham and Women's Hospital, Boston, MA.
- <sup>86</sup>Division of Endocrinology and Metabolism, Department of Internal Medicine, Taichung Veterans General Hospital, Taichung, Taiwan.
- <sup>87</sup>Departments of Medicine, Pharmacology, and Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN.
- <sup>88</sup>Human Genetics Center, University of Texas Health Science Center at Houston and Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX.
- <sup>89</sup>Department of Medicine, Brigham and Women's Hospital, Boston, MA.
- <sup>90</sup>Department of Epidemiology, Harvard School of Public Health, Boston, MA.
- <sup>91</sup>Department of Medicine, Weill Cornell Medical School, New York, NY.
- <sup>92</sup>Institute for Translational Genomics and Population Sciences, Departments of Pediatrics and Medicine, Los Angeles Biomedical Research Institute at Harbor-UCLA Medical Center, Torrance, CA.
- <sup>93</sup>Joseph J. Zilber School of Public Health, University of Wisconsin-Milwaukee, Milwaukee, WI.
- <sup>94</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle WA.
- <sup>95</sup>Department of Biostatistics, School of Public Health, Yale University, New Haven, CT.

<sup>96</sup>Computational Biology and Bioinformatics Program, Yale University, New Haven, CT.

<sup>97</sup>Department of Public Health Sciences, Center for Public Health Genomics, University of Virginia, Charlottesville, VA.

<sup>98</sup>Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS.

<sup>99</sup>Department of Cardiology, Beth Israel Deaconess Medical Center, Boston, MA.

<sup>100</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI.

<sup>101</sup>Department of Biology, MIT, Cambridge, MA, USA.

<sup>102</sup>Department of Systems Biology, Harvard Medical School, Boston, MA, USA.

<sup>103</sup>Harvard Society of Fellows, Harvard University, Cambridge, MA.

<sup>104</sup>Howard Hughes Medical Institute, Boston, MA.

<sup>105</sup>Department of Pathology, Stanford University School of Medicine, Stanford, CA.

<sup>106</sup>Regeneron Pharmaceuticals, Tarrytown, NY.

<sup>107</sup>Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA.

<sup>108</sup>Verve Therapeutics, Cambridge, MA.

## ACKNOWLEDGEMENTS

Investigators who conducted this research report individual research support from R35 HL135818 (S. Redline), P01 HL132825 (S. Weiss), R01HL091357 and R01HL055673 (D. Arnett), W81XWH-17-1-0597 (D. Schwartz), K01 HL135405 (B. Cade), P01 HL132825 (J Lasky-Su), K01HL136700 (S. Aslibekyan), R01HL113323 (J Curran), R01HL1333040 (D. Weeks), 1R01HL138737 (D. Darbar), P01 HL132825 (P. Kachroo), T32 HL129982 (L. Raffield), R01HL113323 (J. Blangero), HHS-N268201800002I (T. Blackwell and A Smith), U54GM115428 (J. Wilson), R01HL148565 and R01HL148050 (P. Natarajan), F30-HL149180 (S. M. Zekavat), R01HL139731 and AHA-18SFRN34250007 (S. Lubitz). Claudia Adams Barr Program for Innovative Cancer Research (V. Sankaran), R01HL142711, MGH Hassenfeld Scholar Award (P. Natarajan), Fondation Leducq TNE-18CVD04 (A. Bick, B. Ebert, S. Jaiswal, P. Natarajan, S. Kathiresan), Burroughs Wellcome Foundation and Ludwig Cancer Center (S. Jaiswal), UMI-HG008895 (S. Kathiresan).

Whole genome sequencing (WGS) for the Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). Centralized read mapping and genotype calling, along with variant quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Phenotype harmonization, data management, sample-identity QC, and general study coordination were provided by the TOPMed Data Coordinating Center (3R01HL-120393-02S1; contract HHSN268201800001I). We gratefully acknowledge the studies and participants who provided biological samples and data for TOPMed.

### COMPETING INTERESTS

B. Psaty serves on the Steering Committee of the Yale Open Data Access Project funded by Johnson & Johnson. E. Silverman and M. Cho received grant support from GlaxoSmithKlein and Bayer. S. Weiss receive royalties from UpToDate. S. Aslibekyan reports Employment/equity in 23andMe, Inc. B. Freedman is a consultant for RenalytixAI and AstraZeneca Pharmaceuticals. MEM reports funding from Regeneron Pharmaceuticals Inc. unrelated to this project. M. Cho has received grant support from GlaxoSmithKlein and Bayer and consulting or speaking fees from AstraZeneca and Illumina. J. Floyd has consulted for Shionogi Inc. B. Levy is a co-founder of Nocien Therapeutics; receives grant support from Pieris Pharmaceuticals, Sanofi and Samsung Research America; and has served as a consultant for Bayer, Entrinsic Health, Gossamer Bio, NControl, Novartis, Teva and Thetis Pharmaceuticals. E. Lander serves on the board of directors for Codiak BioSciences and serves on the scientific

advisory board of F-Prime Capital Partners and Third Rock Ventures. B. Ebert reports grant support from Celgene and Deerfield. P. Ellinor has received grant support from Bayer AG and has served on advisory boards or consulted for Bayer AG, Quest Diagnostics, MyoKardia and Novartis. G. Abecasis is an employee of Regeneron Pharmaceuticals and owns stock and stock options for Regeneron Pharmaceuticals. S. Jaiswal is a scientific advisor to Grail. S. Lubitz receives sponsored research support from Bristol Myers Squibb, Pfizer, Bayer AG, Boehringer Ingelheim, and Fitbit, and has consulted for Bristol Myers Squibb, Pfizer and Bayer AG, and participates in a research collaboration with IBM. P. Natarajan reports grants support from Amgen, Apple, and Boston Scientific, and is a scientific advisor to Apple. S. Kathiresan is an employee of Verve Therapeutics, and holds equity in Verve Therapeutics, Maze Therapeutics, Catabasis, and San Therapeutics. He is a member of the scientific advisory boards for Regeneron Genetics Center and Corvidia Therapeutics; he has served as a consultant for Acceleron, Eli Lilly, Novartis, Merck, Novo Nordisk, Novo Ventures, Ionis, Alnylam, Aegerion, Haug Partners, Noble Insights, Leerink Partners, Bayer Healthcare, Illumina, Color Genomics, MedGenome, Quest, and Medscape. The remainder of the authors report that they have nothing to disclose. The views expressed in this manuscript are those of the authors and do not necessarily represent the views of the National Heart, Lung, and Blood Institute; the National Institutes of Health; or the U.S. Department of Health and Human Services.

## Appendix: NHLBI Trans-Omics for Precision Medicine Consortium

Namiko Abe<sup>109</sup>, Christine Albert<sup>89</sup>, Laura Almasy<sup>110</sup>, Alvaro Alonso<sup>111</sup>, Seth Ament<sup>112</sup>, Peter Anderson<sup>113</sup>, Pramod Anugu<sup>114</sup>, Deborah Applebaum-Bowden<sup>115</sup>, Dan Arking<sup>116</sup>, Allison Ashley-Koch<sup>117</sup>, Stella Aslibekyan<sup>118</sup>, Tim Assimes<sup>119</sup>, Dimitrios Avramopoulos<sup>116</sup>, John Barnard<sup>120</sup>, R. Graham Barr<sup>121</sup>, Emily Barron-Casella<sup>116</sup>, Lucas Barwick<sup>122</sup>, Terri Beaty<sup>116</sup>, Gerald Beck<sup>120</sup>, Diane Becker<sup>59</sup>, Rebecca Beer<sup>115</sup>, Amber Beitelshes<sup>112</sup>, Emelia Benjamin<sup>31</sup>, Panagiotis Benos<sup>123</sup>, Marcos Bezerra<sup>124</sup>, Larry Bielak<sup>4</sup>, Russell Bowler<sup>80</sup>, Jennifer Brody<sup>113</sup>, Ulrich Broeckel<sup>125</sup>, Karen Bunting<sup>109</sup>, Carlos Bustamante<sup>119</sup>, Jonathan Cardwell<sup>126</sup>, Vincent Carey<sup>89</sup>, Cara Carty<sup>127</sup>, Richard Casaburi<sup>128</sup>, James Casella<sup>116</sup>, Peter Castaldi<sup>88</sup>, Mark Chaffin<sup>2</sup>, Christy Chang<sup>112</sup>, Yi-Cheng Chang<sup>129</sup>, Daniel Chasman<sup>89</sup>, Sameer Chavan<sup>126</sup>, Bo-Juen Chen<sup>109</sup>, Wei-Min Chen<sup>130</sup>, Seung Hoan Choi<sup>2</sup>, Lee-Ming Chuang<sup>129</sup>, Mina Chung<sup>120</sup>, Ren-Hua Chung<sup>131</sup>, Clary Clish<sup>2</sup>, Suzy Comhair<sup>12</sup>, Elaine Cornell<sup>69</sup>, Carolyn Crandall<sup>128</sup>, James Crapo<sup>132</sup>, Jeffrey Curtis<sup>4</sup>, Coleen Damcott<sup>112</sup>, Sayantan Das<sup>4</sup>, Sean David<sup>132</sup>, Colleen Davis<sup>113</sup>, Michael DeBaun<sup>133</sup>, Ranjan Deka<sup>134</sup>, Dawn DeMeo<sup>137</sup>, Scott Devine<sup>112</sup>, Qing Duan<sup>135</sup>, Ravi Duggirala<sup>136</sup>, Susan Dutcher<sup>137</sup>, Charles Eaton<sup>138</sup>, Lynette Ekunwe<sup>114</sup>, Adel El Boueiz<sup>29</sup>, Leslie Emery<sup>113</sup>, Serpil Erzurum<sup>120</sup>, Charles Farber<sup>45</sup>, Matthew Flickinger<sup>4</sup>, Nora Franceschini<sup>135</sup>, Chris Frazer<sup>26</sup>, Mao Fu<sup>112</sup>, Stephanie M. Fullerton<sup>113</sup>, Lucinda Fulton<sup>137</sup>, Stacey Gabriel<sup>4</sup>, Weiniu Gan<sup>115</sup>, Shanshan Gao<sup>126</sup>, Yan Gao<sup>114</sup>, Margery Gass<sup>37</sup>, Bruce Gelb<sup>139</sup>, Xiaoqi (Priscilla) Geng<sup>4</sup>, Mark Geraci<sup>140</sup>, Soren Germer<sup>109</sup>, Robert Gerszten<sup>3</sup>, Auyon Ghosh<sup>88</sup>, Richard Gibbs<sup>141</sup>, Chris Gignoux<sup>19</sup>, Mark Gladwin<sup>123</sup>, David Glahn<sup>8</sup>, Da-Wei Gong<sup>112</sup>, Harald Goring<sup>136</sup>, Sharon Graw<sup>126</sup>, Daniel Grine<sup>126</sup>, C. Charles Gu<sup>137</sup>, Yue Guan<sup>112</sup>, Namrata Gupta<sup>3</sup>, Jeff Haessler<sup>126</sup>, Michael Hall<sup>114</sup>, Daniel Harris<sup>112</sup>, Nicola L. Hawley<sup>142</sup>, Ben Heavner<sup>12</sup>, Ryan Hernandez<sup>143</sup>, David Herrington<sup>144</sup>, Craig Hersh<sup>29</sup>, Bertha Hidalgo<sup>61</sup>, Brian Hobbs<sup>88</sup>, John Hokanson<sup>126</sup>, Elliott Hong<sup>112</sup>, Karin Hoth<sup>145</sup>, Chao (Agnes) Hsiung<sup>131</sup>, Yi-Jen Hung<sup>146</sup>, Haley Huston<sup>44</sup>, Chii Min Hwu<sup>147</sup>, Rebecca Jackson<sup>148</sup>, Deepti Jain<sup>12</sup>, Cashell Jaquish<sup>115</sup>, Min A Jhun<sup>4</sup>, Craig Johnson<sup>113</sup>, Rich Johnston<sup>111</sup>, Kimberly Jones<sup>116</sup>, Hyun Min Kang<sup>12</sup>, Shannon Kelly<sup>59</sup>, Michael Kessler<sup>112</sup>, Alyna Khan<sup>113</sup>, Wonji Kim<sup>29</sup>, Greg Kinney<sup>126</sup>, Holly Kramer<sup>149</sup>, Christoph Lange<sup>150</sup>, Meryl LeBoff<sup>47</sup>, Seunggeun Shawn Lee<sup>4</sup>, Wen-Jane Lee<sup>147</sup>, Jonathon LeFaive<sup>4</sup>, David Levine<sup>113</sup>, Joshua Lewis<sup>112</sup>, Xiaohui Li<sup>151</sup>, Yun Li<sup>135</sup>, Henry Lin<sup>151</sup>, Honghuang Lin<sup>152</sup>, Keng Han Lin<sup>4</sup>, Xihong Lin<sup>150</sup>, Simin Liu<sup>153</sup>, Yu Liu<sup>154</sup>, Kathryn Lunetta<sup>152</sup>, James Luo<sup>115</sup>, Michael Mahaney<sup>136</sup>, Barry Make<sup>116</sup>, Ani Manichaikul<sup>97</sup>, Lauren Margolin<sup>2</sup>, Lisa Martin<sup>110</sup>, Susan Mathai<sup>126</sup>, Susanne May<sup>12</sup>,

Patrick McArdle<sup>112</sup>, Merry-Lynn McDonald<sup>118</sup>, Sean McFarland<sup>2,3,5</sup>, Daniel McGoldrick<sup>156</sup>, Caitlin McHugh<sup>12</sup>, Hao Mei<sup>114</sup>, Luisa Mestroni<sup>126</sup>, Julie Mikulla<sup>115</sup>, Nancy Min<sup>114</sup>, Mollie Minear<sup>115</sup>, Ryan L Minster<sup>123</sup>, Matt Moll<sup>85</sup>, Courtney Montgomery<sup>157</sup>, Solomon Musani<sup>41</sup>, Stanford Mwasongwe<sup>114</sup>, Josyf C Mychaleckyj<sup>130</sup>, Girish Nadkarni<sup>158</sup>, Rakhi Naik<sup>116</sup>, Take Naseri<sup>159</sup>, Sergei Nekhai<sup>160</sup>, Sarah C. Nelson<sup>12</sup>, Bonnie Neltner<sup>126</sup>, Deborah Nickerson<sup>156</sup>, Jeff O'Connell<sup>112</sup>, Tim O'Connor<sup>112</sup>, Heather Ochs-Balcom<sup>161</sup>, David Paik<sup>154</sup>, James Pankow<sup>162</sup>, George Papanicolaou<sup>115</sup>, Afshin Parsa<sup>112</sup>, Marco Perez<sup>119</sup>, James Perry<sup>112</sup>, Ulrike Peters<sup>94</sup>, Patricia Peyser<sup>4</sup>, Lawrence S Phillips<sup>111</sup>, Toni Pollin<sup>112</sup>, Wendy Post<sup>116</sup>, Julia Powers Becker<sup>126</sup>, Meher Preethi Boorgula<sup>126</sup>, Michael Preuss<sup>20</sup>, Pankaj Qasba<sup>115</sup>, Dandi Qiao<sup>29</sup>, Zhaohui Qin<sup>111</sup>, Laura Rasmussen-Torvik<sup>72</sup>, Aakrosh Ratan<sup>130</sup>, Robert Reed<sup>112</sup>, Elizabeth Regan<sup>132</sup>, Muagututi'a Sefuiva Reupena<sup>134</sup>, Ken Rice<sup>12</sup>, Carolina Roselli<sup>2</sup>, Ingo Ruczinski<sup>116</sup>, Pamela Russell<sup>126</sup>, Sarah Ruuska<sup>44</sup>, Kathleen Ryan<sup>112</sup>, Danish Saleheen<sup>121</sup>, Shabnam Salimi<sup>112</sup>, Steven Salzberg<sup>116</sup>, Kevin Sandow<sup>151</sup>, Christopher Scheller<sup>4</sup>, Ellen Schmidt<sup>4</sup>, Karen Schwander<sup>137</sup>, Frank Sciorba<sup>123</sup>, Christine Seidman<sup>2,3,15,104</sup>, Jonathan Seidman<sup>2,3</sup>, Vivien Sheehan<sup>164</sup>, Stephanie L. Sherman<sup>165</sup>, Amol Shetty<sup>112</sup>, Aniket Shetty<sup>126</sup>, Brian Silver<sup>166</sup>, Josh Smith<sup>113</sup>, Tanja Smith<sup>109</sup>, Sylvia Smoller<sup>36</sup>, Beverly Snively<sup>144</sup>, Michael Snyder<sup>119</sup>, Tamar Sofer<sup>57</sup>, Nona Sotoodehnia<sup>113</sup>, Adrienne M. Stilp<sup>113</sup>, Garrett Storm<sup>126</sup>, Elizabeth Streeten<sup>112</sup>, Jessica Lasky Su<sup>29</sup>, Yun Ju Sung<sup>137</sup>, Jody Sylvia<sup>47</sup>, Adam Szpiro<sup>113</sup>, Carole Sztalryd<sup>112</sup>, Daniel Taliun<sup>4</sup>, Hua Tang<sup>119</sup>, Matthew Taylor<sup>126</sup>, Simeon Taylor<sup>112</sup>, Marilyn Telen<sup>117</sup>, Timothy A. Thornton<sup>113</sup>, Machiko Threlkeld<sup>156</sup>, Lesley Tinker<sup>126</sup>, David Tirschwell<sup>113</sup>, Sarah Tishkoff<sup>167</sup>, Hemant Tiwari<sup>62</sup>, Catherine Tong<sup>12</sup>, Michael Tsai<sup>162</sup>, Dhananjay Vaidya<sup>116</sup>, David Van Den Berg<sup>168</sup>, Peter VandeHaar<sup>4</sup>, Scott Vrieze<sup>169</sup>, Tarik Walker<sup>126</sup>, Robert Wallace<sup>145</sup>, Avram Walts<sup>126</sup>, Heming Wang<sup>23</sup>, Karol Watson<sup>128</sup>, Bruce Weir<sup>113</sup>, Lu-Chen Weng<sup>2,64</sup>, Jennifer Wessel<sup>170</sup>, Cristen Willer<sup>171</sup>, Kayleen Williams<sup>113</sup>, Carla Wilson<sup>132</sup>, Joseph Wu<sup>154</sup>, Huichun Xu<sup>112</sup>, Lisa Yanek<sup>116</sup>, Rongze Yang<sup>112</sup>, Norann Zaghoul<sup>112</sup>, Yingze Zhang<sup>172</sup>, Snow Xueyan Zhao<sup>32</sup>, Wei Zhao<sup>2</sup>, Degui Zhi<sup>173</sup>, Xiang Zhou<sup>4</sup>, Xiaofeng Zhu<sup>174</sup>, Michael Zody<sup>109</sup>, Sebastian Zoellner<sup>4</sup>

<sup>109</sup>New York Genome Center, New York, NY. <sup>110</sup>Children's Hospital of Philadelphia, University of Pennsylvania, Philadelphia, PA. <sup>111</sup>Emory University, Atlanta, GA. <sup>112</sup>University of Maryland, Baltimore, MD. <sup>113</sup>University of Washington, Seattle, WA. <sup>114</sup>University of Mississippi, Jackson, MS. <sup>115</sup>National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, MD. <sup>116</sup>Johns Hopkins University, Baltimore, MD. <sup>117</sup>Duke University, Durham, NC. <sup>118</sup>University of Alabama, Birmingham, AL. <sup>119</sup>Stanford University, Stanford, CA. <sup>120</sup>Cleveland Clinic, Cleveland, OH. <sup>121</sup>Columbia University, New York, NY. <sup>122</sup>The Emmes Corporation, LTRC, Rockville, MD. <sup>123</sup>University of Pittsburgh, Pittsburgh, PA. <sup>124</sup>Fundação de Hematologia e Hemoterapia de Pernambuco - HEMOPE, Recife, Brazil. <sup>125</sup>Medical College of Wisconsin, Milwaukee, WI. <sup>126</sup>University of Colorado Anschutz Medical Campus, Aurora, CO. <sup>127</sup>Washington State University, Seattle, WA. <sup>128</sup>University of California, Los Angeles, Los Angeles, CA. <sup>129</sup>National Taiwan University Hospital, Taipei, Taiwan. <sup>130</sup>University of Virginia, Charlottesville, VA. <sup>131</sup>National Health Research Institute, Zhunan, Taiwan. <sup>132</sup>National Jewish Health, Denver, CO. <sup>132</sup>University of Chicago, Chicago, IL. <sup>133</sup>Vanderbilt University, Nashville, TN. <sup>134</sup>University of Cincinnati, Cincinnati, OH. <sup>135</sup>University of North Carolina, Chapel Hill, NC. <sup>136</sup>University

of Texas Rio Grande Valley School of Medicine, Edinburg, TX. <sup>137</sup>Washington University in St Louis, St Louis, MO. <sup>138</sup>Brown University, Providence, RI. <sup>139</sup>Icahn School of Medicine at Mount Sinai, New York, NY. <sup>140</sup>Indiana University, Medicine, Indianapolis, IN. <sup>141</sup>Baylor College of Medicine Human Genome Sequencing Center, Houston, TX. <sup>142</sup>Department of Chronic Disease Epidemiology, Yale University, New Haven, CT. <sup>143</sup>McGill University, Montreal, QC, Canada. <sup>144</sup>Wake Forest Baptist Health, Winston-Salem, NC. <sup>145</sup>University of Iowa, Iowa City, IA. <sup>146</sup>Tri-Service General Hospital National Defense Medical Center, Taipei, Taiwan. <sup>147</sup>Taichung Veterans General Hospital Taiwan, Taichung City, Taiwan. <sup>148</sup>Division of Endocrinology, Diabetes and Metabolism, Department of Medicine, Ohio State University Wexner Medical Center, Columbus, OH. <sup>149</sup>Loyola University, Public Health Sciences, Maywood, IL. <sup>150</sup>Department of Biostatistics, Harvard School of Public Health, Boston, MA. <sup>151</sup>Lundquist Institute, Torrance, CA. <sup>152</sup>Boston University, Boston, MA. <sup>153</sup>Department of Epidemiology, Brown University, Providence, RI. <sup>154</sup>Cardiovascular Institute, Stanford University, Palo Alto, CA. <sup>155</sup>George Washington University, Washington, DC. <sup>156</sup>Department of Genome Sciences, University of Washington, Seattle, WA. <sup>157</sup>Oklahoma Medical Research Foundation, Genes and Human Disease, Oklahoma City, OK. <sup>158</sup>Division of Nephrology, Icahn School of Medicine at Mount Sinai, New York, NY. <sup>159</sup>Ministry of Health, Government of Samoa, Apia. <sup>160</sup>Howard University, Washington, DC. <sup>161</sup>University at Buffalo, Buffalo, NY. <sup>162</sup>University of Minnesota, Minneapolis, MN. <sup>163</sup>Lutia i Puava ae Mapu i Fagalele, Apia, Samoa. <sup>164</sup>Department of Pediatrics, Baylor College of Medicine, Houston, TX. <sup>165</sup>Department of Human Genetics, Emory University, Atlanta, GA. <sup>166</sup>UMass Memorial Medical Center, Worcester, MA. <sup>167</sup>Department of Genetics, University of Pennsylvania, Philadelphia, PA. <sup>168</sup>USC Methylation Characterization Center, University of Southern California, Los Angeles, CA. <sup>169</sup>Department of Psychology, University of Minnesota, Minneapolis, MN. <sup>170</sup>Department of Epidemiology, Indiana University, Indianapolis, IN. <sup>171</sup>Department of Medicine, University of Michigan, Ann Arbor, MI. <sup>172</sup>Department of Medicine, University of Pittsburgh, Pittsburgh, PA. <sup>173</sup>Center for Precision Health, School of Biomedical Informatics, University of Texas Health at Houston, Houston, TX. <sup>174</sup>Department of Population and Quantitative Health Sciences, Case Western Reserve University, Cleveland, OH.

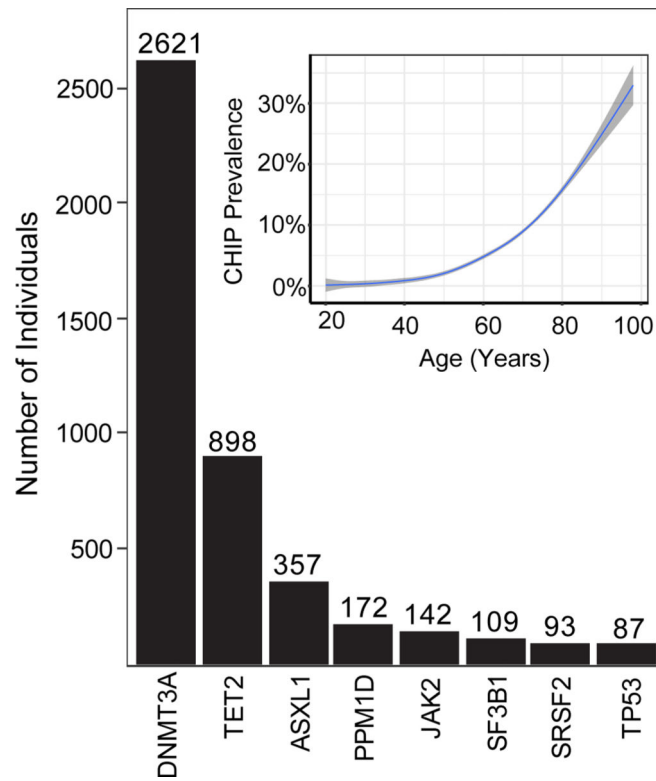
## REFERENCES

1. Kennedy BK et al. Geroscience: linking aging to chronic disease. *Cell* 159, 709–13 (2014). [PubMed: 25417146]
2. Jaiswal S. et al. Age-Related Clonal Hematopoiesis Associated with Adverse Outcomes. *New England Journal of Medicine* 371, 2488–2498 (2014).
3. Genovese G. et al. Clonal Hematopoiesis and Blood-Cancer Risk Inferred from Blood DNA Sequence. *New England Journal of Medicine* 371, 2477–2487 (2014).
4. Xie M. et al. Age-related mutations associated with clonal hematopoietic expansion and malignancies. *Nature Medicine* 20, 1472–1478 (2014).
5. Jaiswal S. et al. Clonal Hematopoiesis and Risk of Atherosclerotic Cardiovascular Disease. *New England Journal of Medicine* 377, 111–121 (2017).
6. Steensma DP et al. Clonal hematopoiesis of indeterminate potential and its distinction from myelodysplastic syndromes. *Blood* 126, 9–16 (2015). [PubMed: 25931582]



7. Taliun D. et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *bioRxiv* 563866 (2019). doi:10.1101/563866
8. Cibulskis K. et al. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature Biotechnology* (2013). doi:10.1038/nbt.2514
9. Loh PR et al. Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature* (2018). doi:10.1038/s41586-018-0321-x
10. Jaiswal S. et al. Age-Related Clonal Hematopoiesis Associated with Adverse Outcomes. *New England Journal of Medicine* 371, 2488–2498 (2014).
11. Patel K. v. et al. Red Cell distribution width and mortality in older adults: A meta-analysis. *Journals of Gerontology - Series A Biological Sciences and Medical Sciences* 65 A, 258–265 (2010).
12. Bick AG et al. Genetic Interleukin 6 Signaling Deficiency Attenuates Cardiovascular Risk in Clonal Hematopoiesis. *Circulation* 141, 124–131 (2020). [PubMed: 31707836]
13. Alexandrov LB et al. Clock-like mutational processes in human somatic cells. *Nature Genetics* 47, 1402–1407 (2015). [PubMed: 26551669]
14. Zink F. et al. Clonal hematopoiesis, with and without candidate driver mutations, is common in the elderly. *Blood* 130, 742–752 (2017). [PubMed: 28483762]
15. Bowman RL, Busque L & Levine RL Cell Stem Cell Review Clonal Hematopoiesis and Evolution to Hematopoietic Malignancies. (2018). doi:10.1016/j.stem.2018.01.011
16. Desai P. et al. Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nature Medicine* 24, 1015–1023 (2018).
17. Bojesen SE et al. Multiple independent variants at the TERT locus are associated with telomere length and risks of breast and ovarian cancer. *Nature Genetics* (2013). doi:10.1038/ng.2566
18. Zhou W. et al. Mosaic loss of chromosome Y is associated with common variation near TCL1A. *Nature genetics* 48, 563–8 (2016). [PubMed: 27064253]
19. Hinds DA et al. Germ line variants predispose to both JAK2 V617F clonal hematopoiesis and myeloproliferative neoplasms. *Blood* (2016). doi:10.1182/blood-2015-06-652941
20. Hu Y. et al. A statistical framework for cross-tissue transcriptome-wide association analysis. *Nature Genetics* (2019). doi:10.1038/s41588-019-0345-7
21. Smith BW et al. The aryl hydrocarbon receptor directs hematopoietic progenitor cell expansion and differentiation. *Blood* (2013). doi:10.1182/blood-2012-11-466722
22. Cybulski C. et al. CHEK2 Is a Multiorgan Cancer Susceptibility Gene. *The American Journal of Human Genetics* 75, 1131–1135 (2004). [PubMed: 15492928]
23. Rudd MF, Sellick GS, Webb EL, Catovsky D & Houlston RS Variants in the ATM-BRCA2-CHEK2 axis predispose to chronic lymphocytic leukemia. *Blood* (2006). doi:10.1182/blood-2005-12-5022
24. Huynh M. et al. Hyaluronan and proteoglycan link protein 1 (HAPLN1) activates bortezomib-resistant NF- $\kappa$ B activity and increases drug resistance in multiple myeloma. *The Journal of biological chemistry* 293, 2452–2465 (2018). [PubMed: 29279332]
25. Moran-Crusio K. et al. Tet2 Loss Leads to Increased Hematopoietic Stem Cell Self-Renewal and Myeloid Transformation. *Cancer Cell* (2011). doi:10.1016/j.ccr.2011.06.001
26. Kilpivaara O. et al. A germline JAK2 SNP is associated with predisposition to the development of JAK2V617F-positive myeloproliferative neoplasms. *Nature Genetics* 41, 455–459 (2009). [PubMed: 19287384]
27. Jones A. v. et al. JAK2 haplotype is a major risk factor for the development of myeloproliferative neoplasms. *Nature Genetics* 41, 446–449 (2009). [PubMed: 19287382]
28. Olcaydu D. et al. A common JAK2 haplotype confers susceptibility to myeloproliferative neoplasms. *Nature Genetics* 41, 450–454 (2009). [PubMed: 19287385]
29. Young AL, Challen GA, Birman BM & Druley TE Clonal haematopoiesis harbouring AML-associated mutations is ubiquitous in healthy adults. *Nature Communications* 7, 1–7 (2016).
30. Regier AA et al. Functional equivalence of genome sequencing analysis pipelines enables harmonized variant calling across human genetics projects. *Nature Communications* (2018). doi:10.1038/s41467-018-06159-4

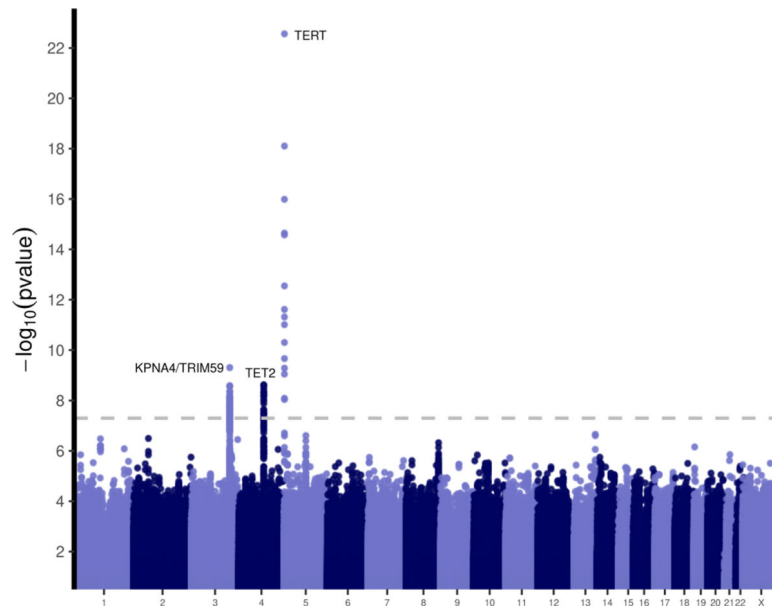
31. Jun G, Wing MK, Abecasis GR & Kang HM An efficient and scalable analysis framework for variant extraction and refinement from population-scale DNA sequence data. *Genome Research* 25, 918–925 (2015). [PubMed: 25883319]
32. Karczewski KJ et al. Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv* 531210 (2019). doi:10.1101/531210
33. Gibson CJ et al. Clonal Hematopoiesis Associated With Adverse Outcomes After Autologous Stem-Cell Transplantation for Lymphoma. *Journal of clinical oncology : official journal of the American Society of Clinical Oncology* 35, 1598–1605 (2017). [PubMed: 28068180]
34. Hiatt JB, Pritchard CC, Salipante SJ, O’Roak BJ & Shendure J Single molecule molecular inversion probes for targeted, high-accuracy detection of low-frequency variation. *Genome Research* 23, 843–854 (2013). [PubMed: 23382536]
35. Pérez Millán MI et al. Next generation sequencing panel based on single molecule molecular inversion probes for detecting genetic variants in children with hypopituitarism. *Molecular Genetics and Genomic Medicine* 6, 514–525 (2018).
36. Li Y, Willer CJ, Ding J, Scheet P & Abecasis GR MaCH: Using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genetic Epidemiology* 34, 816–834 (2010). [PubMed: 21058334]
37. Vattathil S & Scheet P Haplotype-based profiling of subtle allelic imbalance with SNP arrays. *Genome Research* 23, 152–158 (2013). [PubMed: 23028187]
38. Fowler J, San Lucas FA & Scheet P System for Quality-Assured Data Analysis: Flexible, reproducible scientific workflows. *Genetic Epidemiology* 43, 227–237 (2019). [PubMed: 30565316]
39. Natarajan P. et al. Deep-coverage whole genome sequences and blood lipids among 16,324 individuals. *Nature Communications* 9, 3391 (2018).
40. Nik-Zainal S. et al. Landscape of somatic mutations in 560 breast cancer whole-genome sequences. *Nature* 534, 47–54 (2016). [PubMed: 27135926]
41. Blokzijl F, Janssen R, van Boxtel R & Cuppen E MutationalPatterns: Comprehensive genome-wide analysis of mutational processes. *Genome Medicine* 10, 33 (2018). [PubMed: 29695279]
42. Zhou W. et al. Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nature genetics* 50, 1335–1341 (2018). [PubMed: 30104761]
43. Bates D. et al. Package “lme4”: Linear Mixed-Effects Models using “Eigen” and S4. *Journal of Statistical Software* (2015). doi:10.18637/jss.v067.i01
44. Benner C. et al. FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics (Oxford, England)* 32, 1493–1501 (2016).
45. Amemiya HM, Kundaje A & Boyle AP The ENCODE Blacklist: Identification of Problematic Regions of the Genome. *Scientific reports* 9, 9354 (2019). [PubMed: 31249361]
46. Fulco CP et al. Activity-by-contact model of enhancer-promoter regulation from thousands of CRISPR perturbations. *Nature Genetics* doi:10.1038/s41588-019-0538-0
47. Johnson WE, Li C & Rabinovic A Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics (Oxford, England)* 8, 118–27 (2007).
48. Dobin A. et al. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15–21 (2013). [PubMed: 23104886]
49. DeLuca DS et al. RNA-SeQC: RNA-seq metrics for quality control and process optimization. *Bioinformatics (Oxford, England)* 28, 1530–2 (2012).
50. Stranger BE et al. Enhancing GTEx by bridging the gaps between genotype, gene expression, and disease. *Nature Genetics* 49, 1664–1670 (2017). [PubMed: 29019975]
51. Houseman EA et al. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics* (2012). doi:10.1186/1471-2105-13-86
52. Horvath S & Levine AJ HIV-1 infection accelerates age according to the epigenetic clock. *Journal of Infectious Diseases* (2015). doi:10.1093/infdis/jiv277
53. Barfield RT, Kilaru V, Smith AK & Conneely KN CpGassoc: An R function for analysis of DNA methylation microarray data. *Bioinformatics* (2012). doi:10.1093/bioinformatics/bts124



**Fig. 1|. Identifying CHIP in TOPMed Genomes.**

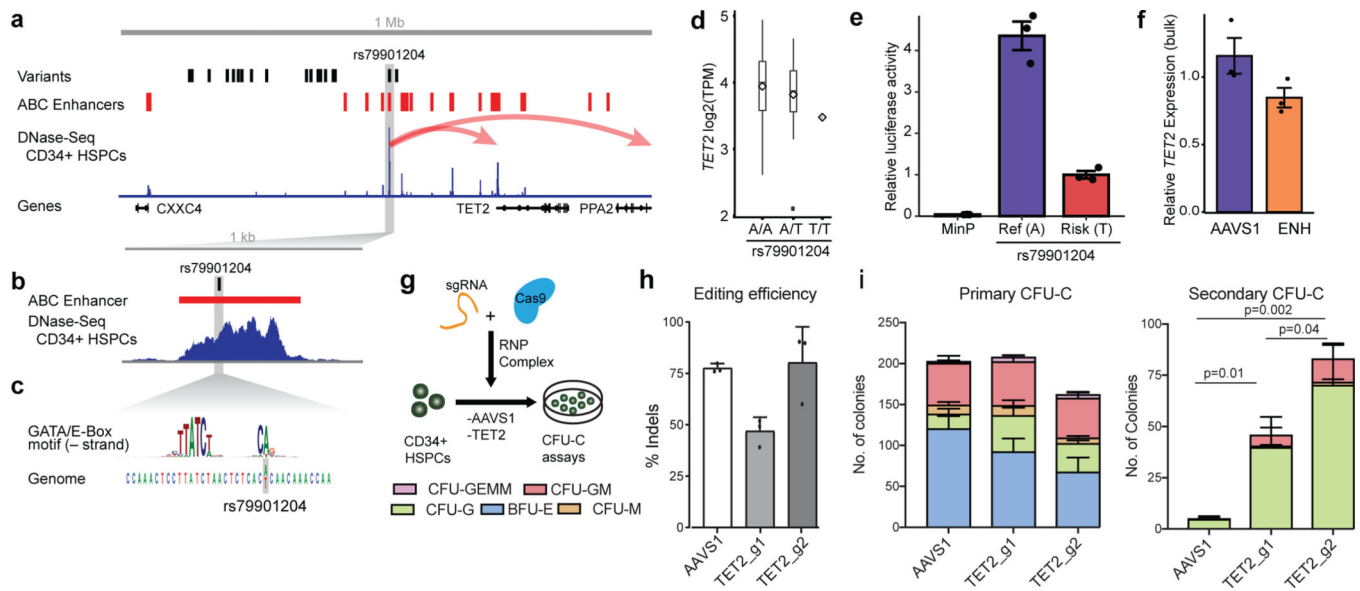
CHIP was identified in 97,631 whole genome sequenced peripheral blood samples through the curation of somatic driver mutations. Counts for 8 most common driver genes plotted.

**inset**, CHIP prevalence increased with age. Center line represents general additive model spline, 95% confidence interval is shaded (n=82,807 individuals; two-sided t-test:  $p < 10^{-300}$ ).



**Fig. 2|. Genetic Determinants of CHIP.**

Single variant genetic association analyses of CHIP identified three genome wide significant loci. Two-sided association testing performed using SAIGE (n=65,405 individuals, see methods)



**Fig. 3]. African ancestry specific *TET2* locus risk variant disrupts hematopoietic stem cell *TET2* enhancer decreasing *TET2* expression and increasing self-renewal.**

**a**, the *TET2* locus with fine-mapped risk variants, Activity-by-Contact (ABC) hematopoietic stem and progenitor cell (HSPC) enhancers, DNase-Seq CD34+ HSPC and RefSeq genes. ABC model predicts that rs79901204 disrupts a *TET2* enhancer resulting in decreased *TET2* expression (see methods). **b**, expanded view of *TET2* enhancer element. **c**, rs79901204 disrupts a GATA motif/E-Box motif. **d**, rs79901204 is associated with decreased *TET2* expression in human peripheral blood RNA-seq ( $N_{A/A}=230$ ,  $N_{A/T}=16$ ,  $N_{T/T}=1$ , two-sided linear mixed model  $p=0.012$ ). TPM, transcripts per million. Boxplot displays median, 25<sup>th</sup> and 75<sup>th</sup> percentiles, mean (diamond symbol) and outlier observations (black dots) **e**, luciferase assay in CD34+ primary cells demonstrates four-fold attenuation of enhancer activity by the rs79901204 T risk allele relative to the A reference allele ( $N=3$ , two-sided t-test  $p=0.007$ ). **f**, deleting the *TET2* enhancer (ENH) in CD34+ primary cells results in decreased *TET2* expression relative to deletion of control locus AAVS1 ( $N=3$ , two-sided t-test,  $p=0.04$ ). **g**, Human HSPCs were electroporated with Cas9 targeting a coding region of *TET2* and AAVS1 (a control locus) and plated for primary and secondary colony-forming assays. **h**, two *TET2* guides had differential editing efficiency. **i**, *TET2* coding disruption leads to expanded secondary colony formation compared to AAVS1 controls ( $N=3$ , two-sided t-test  $p=0.01$ ,  $p=0.002$  for g1 and g2 respectively, with greater expansion identified in the *TET2* guide with greater editing efficiency (two-sided t-test  $p=0.04$ ). Mean and standard deviation of number of each colony type plotted. CFU-M, colony forming unit-macrophage; CFU-GM, granulocyte macrophage; CFU-GEMM, granulocyte erythrocyte macrophage megakaryocyte; CFU-G, granulocyte; BFU-E, burst forming unit-erythroid. In **e**, **f**, **h**, points represent independent replicates, mean values and error bars represent standard error are plotted.