**Supplemental information**

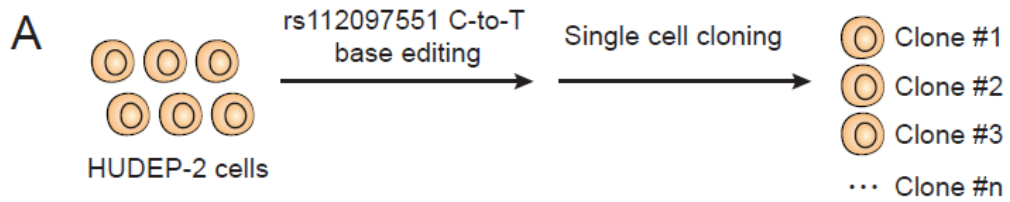# Whole-genome sequencing association analysis

# of quantitative red blood cell phenotypes:

# The NHLBI TOPMed program

Yao Hu, Adrienne M. Stilp, Caitlin P. McHugh, Shuquan Rao, Deepti Jain, Xiuwen Zheng, John Lane, Sébastian Méric de Bellefon, Laura M. Raffield, Ming-Huei Chen, Lisa R. Yanek, Marsha Wheeler, Yao Yao, Chunyan Ren, Jai Broome, Jee-Young Moon, Paul S. de Vries, Brian D. Hobbs, Quan Sun, Praveen Surendran, Jennifer A. Brody, Thomas W. Blackwell, Hélène Choquet, Kathleen Ryan, Ravindranath Duggirala, Nancy Heard-Costa, Zhe Wang, Nathalie Chami, Michael H. Preuss, Nancy Min, Lynette Ekunwe, Leslie A. Lange, Mary Cushman, Nauder Faraday, Joanne E. Curran, Laura Almasy, Kousik Kundu, Albert V. Smith, Stacey Gabriel, Jerome I. Rotter, Myriam Fornage, Donald M. Lloyd-Jones, Ramachandran S. Vasan, Nicholas L. Smith, Kari E. North, Eric Boerwinkle, Lewis C. Becker, Joshua P. Lewis, Goncalo R. Abecasis, Lifang Hou, Jeffrey R. O'Connell, Alanna C. Morrison, Terri H. Beaty, Robert Kaplan, Adolfo Correa, John Blangero, Eric Jorgenson, Bruce M. Psaty, Charles Kooperberg, Russell T. Walton, Benjamin P. Kleinstiver, Hua Tang, Ruth J.F. Loos, Nicole Soranzo, Adam S. Butterworth, Debbie Nickerson, Stephen S. Rich, Braxton D. Mitchell, Andrew D. Johnson, Paul L. Auer, Yun Li, Rasika A. Mathias, Guillaume Lettre, Nathan Pankratz, Cathy C. Laurie, Cecelia A. Laurie, Daniel E. Bauer, Matthew P. Conomos, Alexander P. Reiner, and and the NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium

**Supplemental Figures and Legends**

**Figure S1. Rs112097551 C-to-T base editing and single cell cloning in HUDEP-2 cells. (A) Scheme of rs112097551 C-to-T base editing and FACS-based single cell separation. (B) Efficiency of rs112097551 C-to-T (G-to-A on opposing strand) base editing efficiency in all five clones. Since base editor and sgRNA are constitutively expressed, the frequency of C-to-T conversion may exceed 50% in heterozygous clones due to base editing after single cell cloning.**
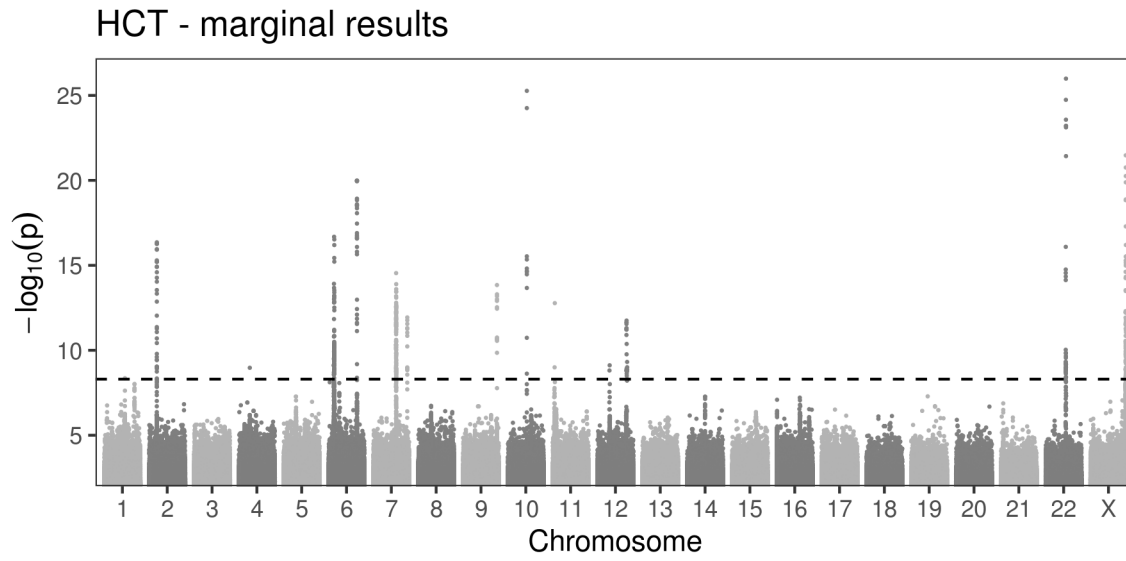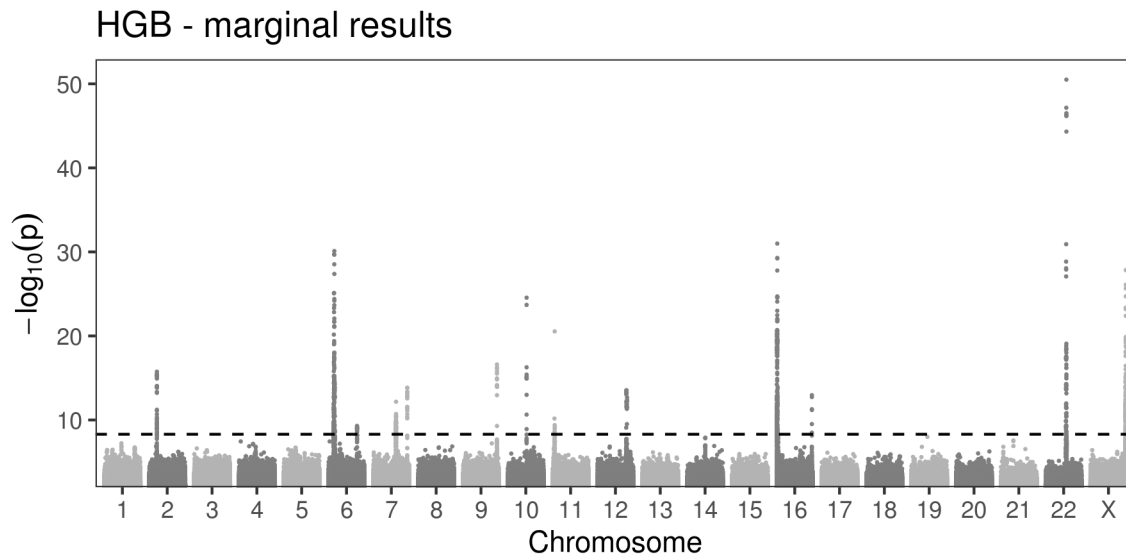
A

rs112097551 C-to-T base editing → Single cell cloning →

HUDEP-2 cells

Clone #1
Clone #2
Clone #3
··· Clone #n

B

| Clone ID | Allele A percentage |
|----------|---------------------|
| Clone #1 | 59% |
| Clone #2 | 61% |
| Clone #3 | 65% |
| Clone #4 | 54% |
| Clone #5 | 58% |

**Figure S2. Manhattan plots of the marginal single-variant analyses in TOPMed. (A) HCT; (B) HGB; (C) MCH; (D) MCHC; (E) MCV; (F) RBC; (G) RDW.**
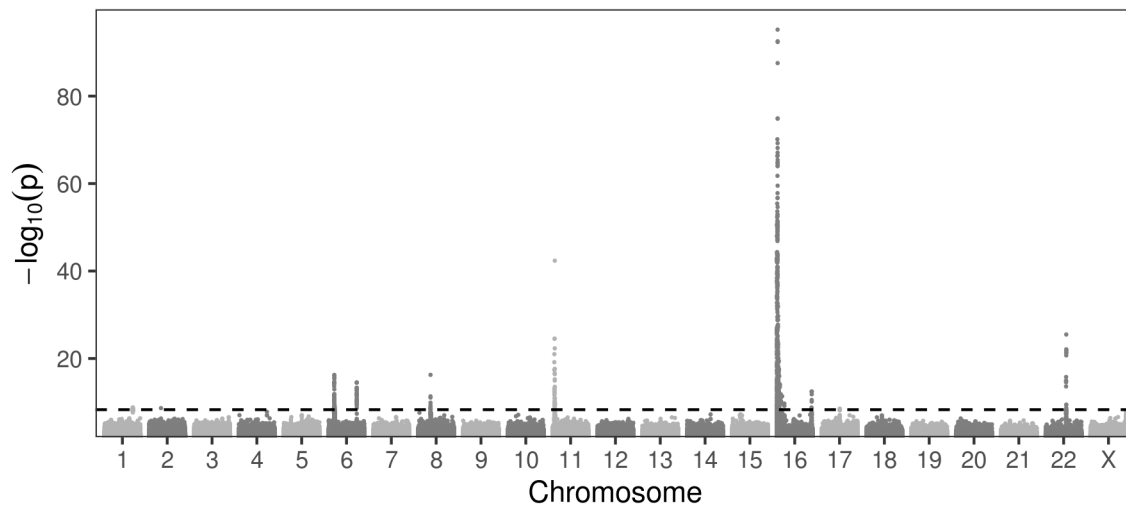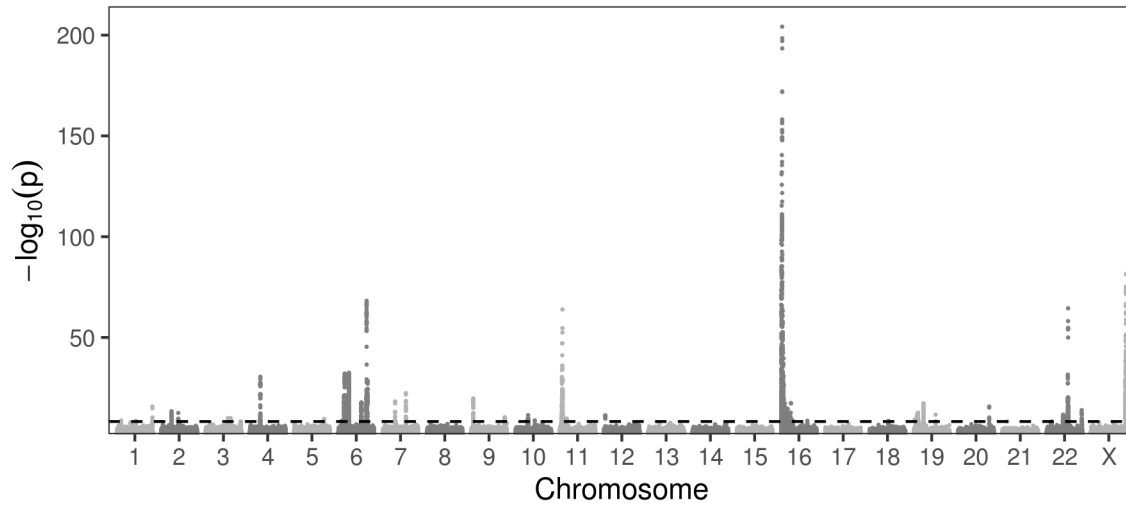
(A)



(B)

(C)

## MCH - marginal results



(D)
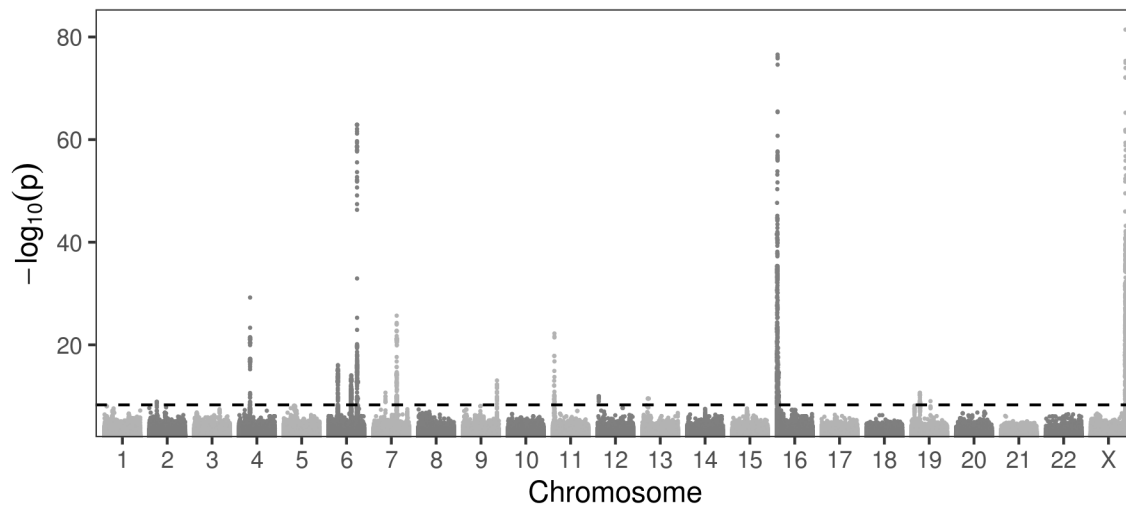
## MCHC - marginal results

(E)

MCV - marginal results


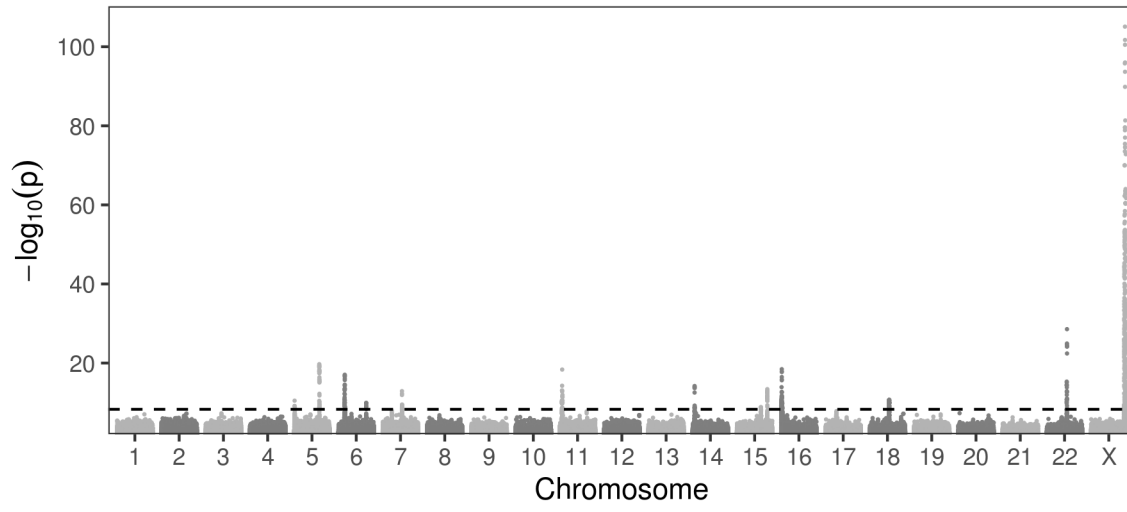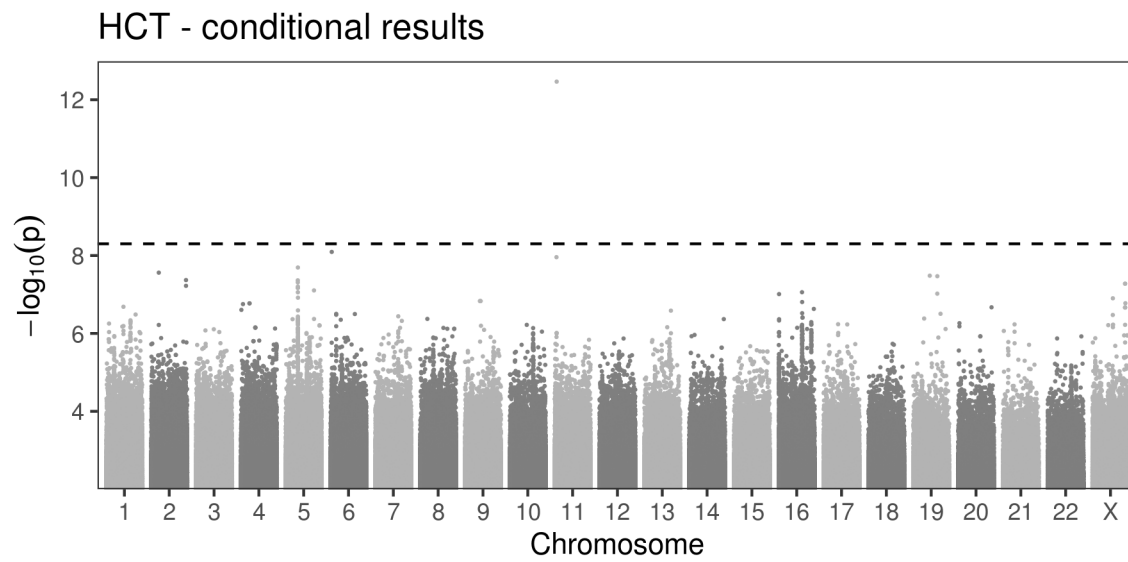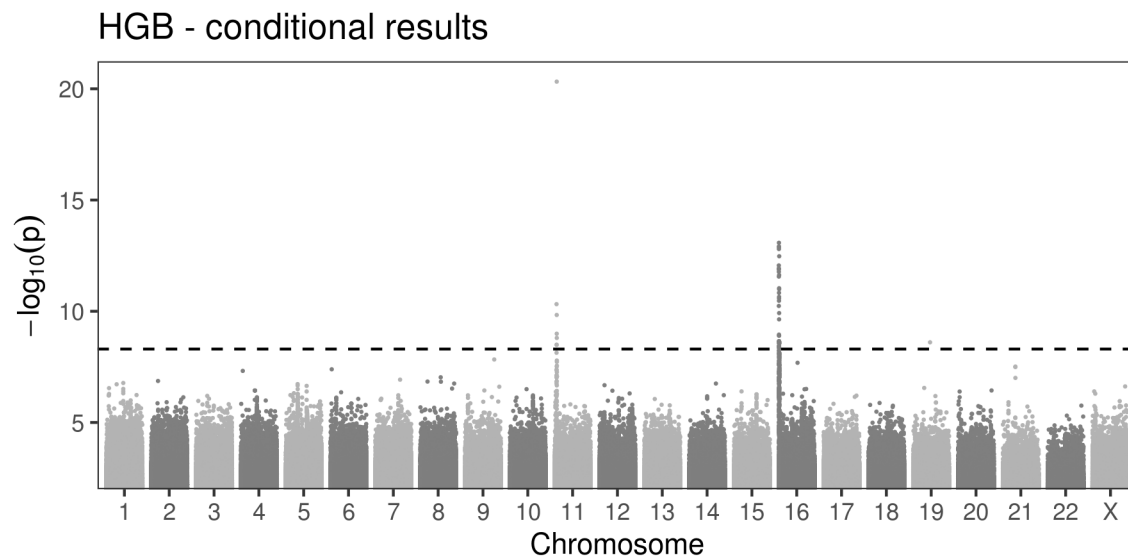
(F)

RBC - marginal results

(G)

RDW - marginal results

**Figure S3. Manhattan plots of the trait-specific conditional single-variant analyses in TOPMed. (A) HCT; (B) HGB; (C) MCH; (D) MCHC; (E) MCV; (F) RBC; (G) RDW.**
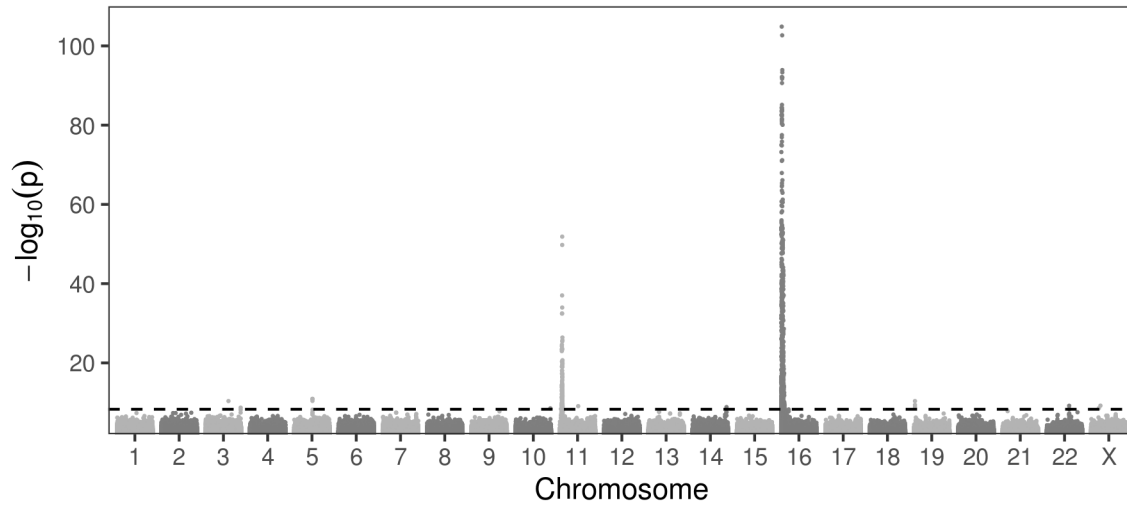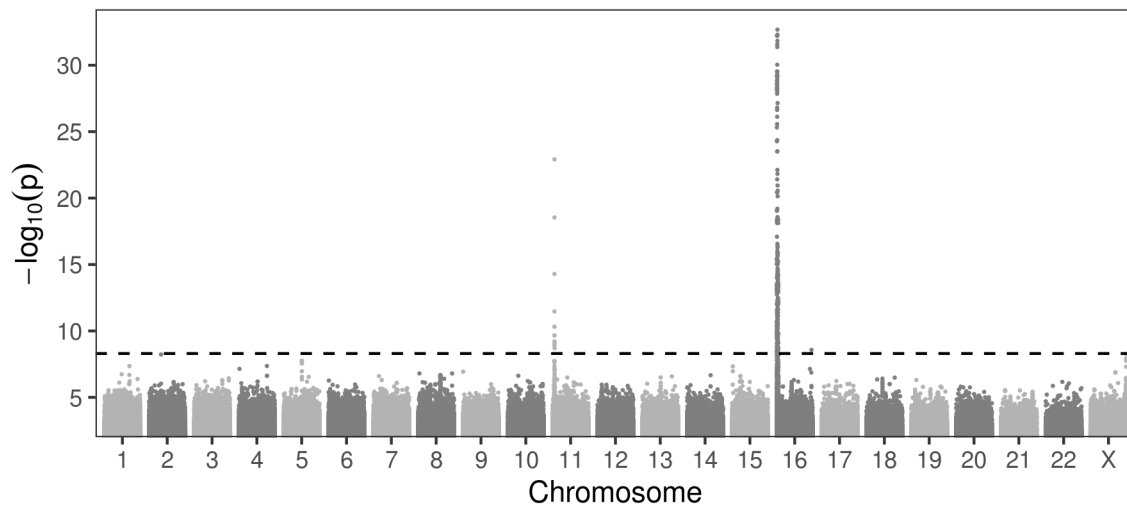
(A)



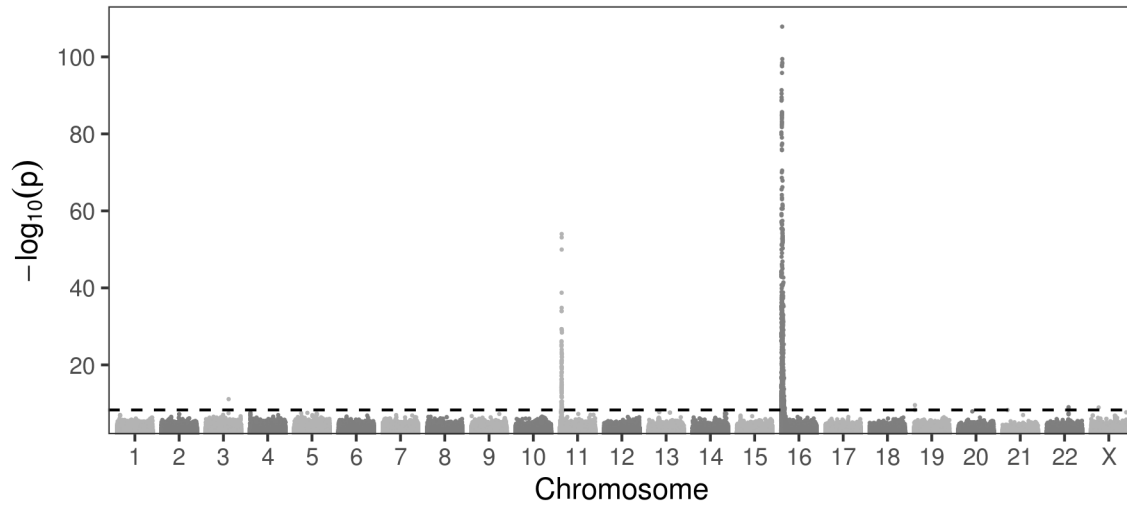(B)

(C)

## MCH - conditional results



(D)

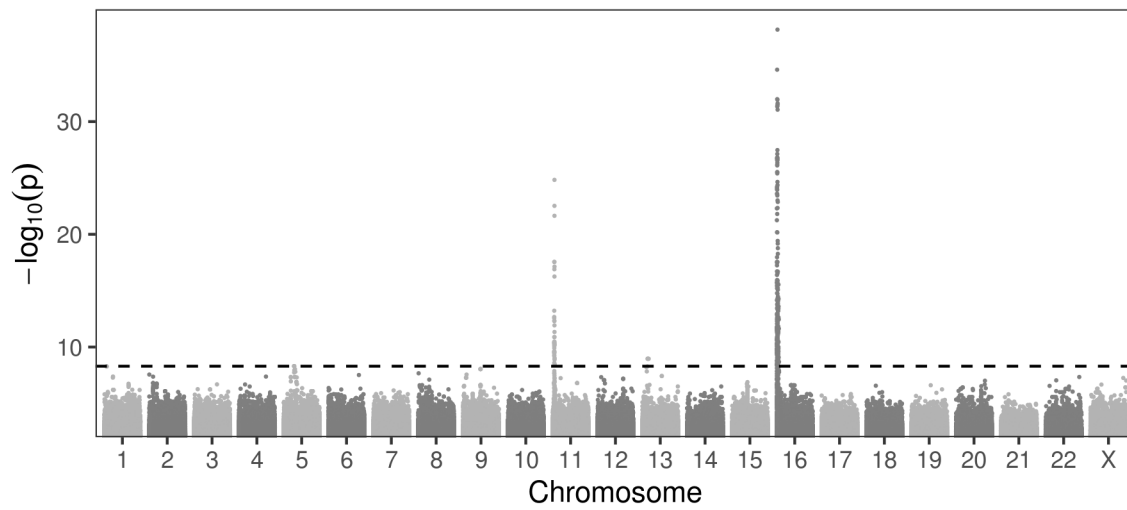## MCHC - conditional results

(E)

## MCV - conditional results



(F)

## RBC - conditional results
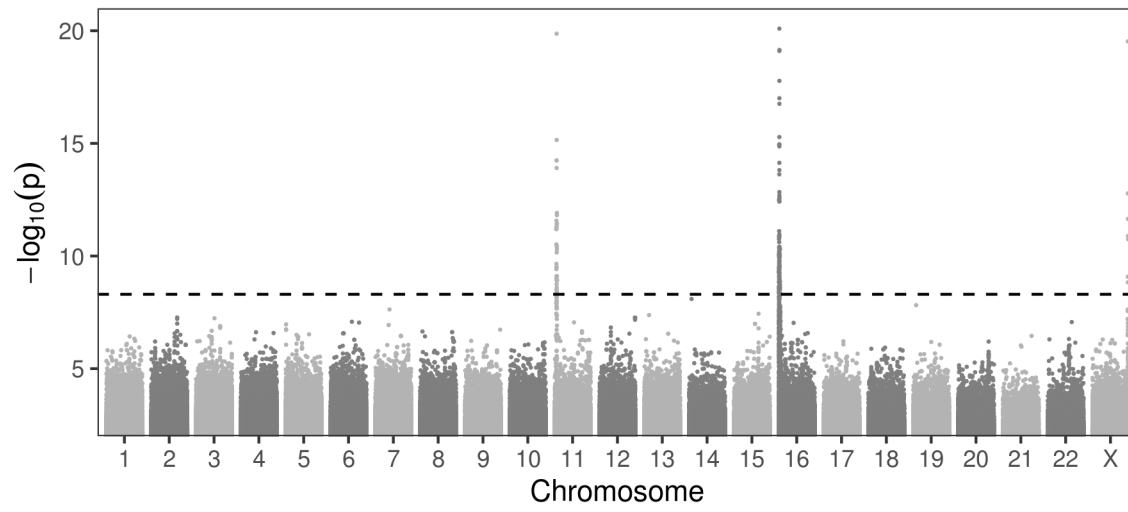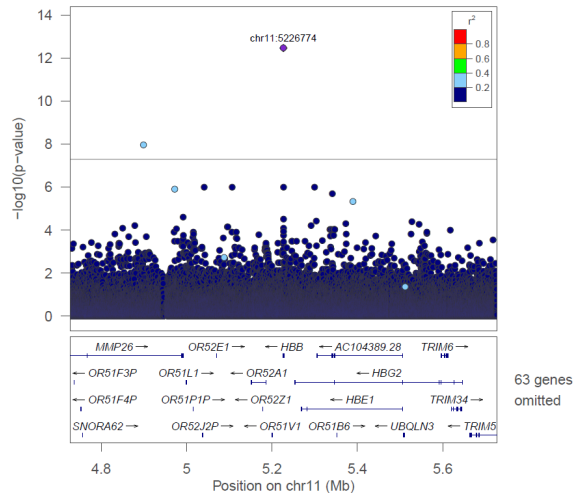
(G)

RDW - conditional results

**Figure S4. Locuszoom plots of the 12 novel variants and conditionally independent variants identified in TOPMed. (A) HCT; (B) HGB; (C) MCH; (D) MCHC; (E) MCV; (F) RBC; (G) RDW**
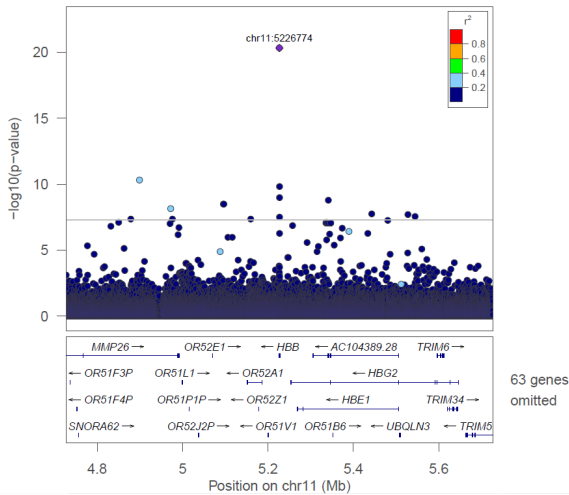
(A)



(B)

(C)

(D)

(E)

(F)

(G)

**Figure S5. Rare variants identified in the aggregated analysis in TOPMed. (A) HCT; (B) HGB; (C) MCH; (D) MCHC; (E) MCV; (F) RBC; (G) RDW.**

(A)

(A1)

(A2)

HCT - HBB - coding2_relaxed

p = 6.36e-10



(A3)

HCT - HBB - coding1_stringent

p = 3.82e-13

(A4)

HCT - HBB - coding1_relaxed

p = 2.79e-09

(A5)

HCT - HBB - coding2_relaxed + noncoding_stringent

p = 6.36e-10

(B)

(B1)



HGB - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 5.15e-08

(B2)

HGB - HBB - coding2_relaxed

p = 9.77e-09

(B3)

HGB - HBB - coding1_stringent

p = 3.52e-11

**(B4)**



HGB - HBB - coding1_relaxed

p = 5.99e-08

**(B5)**



HGB - HBB - coding2_relaxed + noncoding_stringent

p = 9.77e-09

(C)

(C1)

MCH - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 9.07e-11

(C2)



MCH - G6PD - coding2_relaxed

p = 1.23e-06

(C3)



MCH - G6PD - coding1_relaxed

p = 1.12e-06

(C4)

MCH - G6PD - coding2_relaxed + noncoding_relaxed

p = 1.23e-06



(C5)

MCH - G6PD - coding2_relaxed + noncoding_stringent

p = 1.23e-06

(C6)

## MCH - HBA1 - coding2_relaxed

p = 1.82e-09



(C7)

## MCH - HBA1 - coding1_stringent

p = 4.88e-07

## (C8)

### MCH - HBA1 - coding1_relaxed

p = 3.04e-09



## (C9)

### MCH - HBA1 - coding2_relaxed + noncoding_relaxed

p = 1.82e-09

(C10)

MCH - HBA1 - coding2_relaxed + noncoding_stringent

p = 1.82e-09

(C11)

MCH - HBB - coding2_relaxed

p = 1.04e-19

(C12)

MCH - HBB - coding1_stringent

p = 1.07e-14



(C13)

MCH - HBB - coding1_relaxed

p = 5.62e-19

(C14)

MCH - HBB - coding2_relaxed + noncoding_relaxed

p = 7.09e-11



(C15)

MCH - HBB - coding2_relaxed + noncoding_stringent

p = 1.04e-19

(C16)

MCH - TMPRSS6 - coding2_relaxed

p = 9.88e-08



(C17)

MCH - TMPRSS6 - coding1_stringent

p = 4.91e-11

(C18)

MCH - TMPRSS6 - coding1_relaxed

p = 1.32e-11



(C19)

MCH - TMPRSS6 - coding2_relaxed + noncoding_stringent

p = 2.35e-07

(D)

(D1)

MCHC - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 8.76e-10

(D2)

MCHC - HBB - coding2_relaxed

p = 2.94e-12

(D3)

MCHC - HBB - coding1_stringent

p = 2.83e-12

(D4)

MCHC - HBB - coding1_relaxed

p = 2.06e-11



(D5)

MCHC - HBB - coding2_relaxed + noncoding_stringent

p = 2.94e-12

(E)

(E1)



MCV - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 8.6e-14

(E2)

MCV - G6PD - coding2_relaxed

p = 1.63e-13



(E3)

MCV - G6PD - coding1_relaxed

p = 2.88e-13

(E4)

MCV - G6PD - coding2_relaxed + noncoding_relaxed

p = 1.63e-13



(E5)

MCV - G6PD - coding2_relaxed + noncoding_stringent

p = 1.63e-13

(E6)

MCV - HBA1 - coding2_relaxed

p = 2.53e-06



(E7)

MCV - HBB - coding2_relaxed

p = 4.27e-24

(E8)

MCV - HBB - coding1_stringent

p = 3.88e-20



(E9)

MCV - HBB - coding1_relaxed

p = 5.63e-24

(E10)

MCV - HBB - coding2_relaxed + noncoding_relaxed

p = 4.92e-11



(E11)

MCV - HBB - coding2_relaxed + noncoding_stringent

p = 4.27e-24

(E12)

MCV - TFRC - coding1_relaxed

p = 1.57e-06



(E13)

MCV - TMPRSS6 - coding1_stringent

p = 1.62e-08

(E14)



MCV - TMPRSS6 - coding1_relaxed

p = 2.61e-09

(F)

(F1)



RBC - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 1.02e-11

(F2)

RBC - G6PD - coding2_relaxed

p = 3.07e-09



(F3)

RBC - G6PD - coding1_relaxed

p = 1.63e-09

(F4)

RBC - G6PD - coding2_relaxed + noncoding_relaxed

p = 3.07e-09



(F5)

RBC - G6PD - coding2_relaxed + noncoding_stringent

p = 3.07e-09

(F6)

RBC - HBB - coding2_relaxed

p = 1.54e-14

(F7)

RBC - HBB - coding1_stringent

p = 8.83e-14

(F8)



RBC - HBB - coding1_relaxed

p = 6.03e-15

(F9)



RBC - HBB - coding2_relaxed + noncoding_stringent

p = 1.54e-14

(G)

(G1)

RDW - CD36 - coding2_relaxed

p = 2.17e-06

(G2)

RDW - CD36 - coding1_stringent

p = 3.47e-07



(G3)

RDW - AC104389.6 - coding2_relaxed + noncoding_relaxed

p = 3.62e-12

(G4)

RDW - G6PD - coding2_relaxed

p = 2.53e-10

(G5)

RDW - G6PD - coding1_relaxed

p = 2.07e-11

(G6)

RDW - G6PD - coding2_relaxed + noncoding_relaxed

p = 2.53e-10



(G7)

RDW - G6PD - coding2_relaxed + noncoding_stringent

p = 2.53e-10

(G8)

RDW - HBB - coding2_relaxed

p = 1.02e-10

(G9)

RDW - HBB - coding1_stringent

p = 5.58e-15

## (G10)

### RDW - HBB - coding1_relaxed

p = 3.42e-11



## (G11)

### RDW - HBB - coding2_relaxed + noncoding_relaxed

p = 1.58e-07

(G12)

RDW - HBB - coding2_relaxed + noncoding_stringent

p = 1.02e-10

(G13)

RDW - SLC12A7 - coding1_relaxed

p = 2.52e-06

(G14)

RDW - TMPRSS6 - coding1_stringent

p = 1.87e-07

**Supplemental Tables**

**Table S1. Counts of participants by HARE group for each RBC phenotype**

| | Amish | Asian | Black | Central American | Cuban | Dominican | Mexican | Puerto Rican | South American | White |
|---|---|---|---|---|---|---|---|---|---|---|
| **HCT** | 1102 | 654 | 14474 | 708 | 2037 | 2049 | 3556 | 4977 | 708 | 32222 |
| **HGB** | 1102 | 653 | 14454 | 708 | 2037 | 2048 | 3556 | 4974 | 706 | 32223 |
| **MCH** | 1102 | 447 | 11246 | 708 | 2002 | 2049 | 3435 | 4934 | 706 | 19612 |
| **MCHC** | 1102 | 447 | 13112 | 708 | 2002 | 2049 | 3434 | 4934 | 706 | 24154 |
| **MCV** | 1102 | 447 | 12285 | 708 | 2002 | 2049 | 3432 | 4934 | 706 | 21165 |
| **RBC** | 1102 | 384 | 10747 | 682 | 1984 | 1938 | 3413 | 4448 | 654 | 19118 |
| **RDW** | 0 | 447 | 6776 | 662 | 2002 | 1936 | 1898 | 4833 | 647 | 10184 |

HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width.

**Table S2. Pairwise trait correlation (upper triangle) and the number of samples used to calculate the correlations (lower triangle)**

|       | HCT   | HGB   | MCH     | MCHC    | MCV     | RBC      | RDW      |
|-------|-------|-------|---------|---------|---------|----------|----------|
| **HCT**  |       | 0.93211 | 0.21225 | 0.03469 | 0.23772 | 0.74892  | -0.27781 |
| **HGB**  | 62447 |       | 0.35214 | 0.31502 | 0.27268 | 0.70519  | -0.38713 |
| **MCH**  | 46099 | 46083 |         | 0.52655 | 0.87826 | -0.33338 | -0.44973 |
| **MCHC** | 52628 | 52612 | 46109   |         | 0.16304 | -0.05466 | -0.37418 |
| **MCV**  | 48807 | 48791 | 46116   | 48816   |         | -0.3458  | -0.3531  |
| **RBC**  | 44326 | 44309 | 44430   | 44334   | 44340   |          | -0.0827  |
| **RDW**  | 29244 | 29235 | 29350   | 29254   | 29261   | 27572    |          |

HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width.

**Table S3. Basic characteristics of each participating study in TOPMed stratified by race/ethnicity**
See Excel file.

**Table S4. Number of SNVs and indels tested for each RBC trait in TOPMed**

| Trait | Indel | SNV | Total |
|-------|-------|-----|-------|
| RDW | 5356149 | 70775433 | 76131582 |
| RBC | 6497451 | 86154571 | 92652022 |
| MCH | 6637757 | 88043681 | 94681438 |
| MCV | 6834916 | 90671951 | 97506867 |
| MCHC | 7089902 | 94110688 | 101200590 |
| HGB | 7719013 | 102632640 | 110351653 |
| HCT | 7722116 | 102674666 | 110396782 |

HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width.

**Table S5. Previously reported variants and indication of inclusion in the conditional analysis in TOPMed**
See Excel file.

**Table S6. Guide sequence used in the present study**

| Guide | Sequence (5'-3') |
|---|---|
| Guides for CRISPR/Cas9 editing | |
| RUVBL1 | ACTACTTACCAATGGCCCTG |
| Neutral locus | GTAAGCTTAAAACATTAGTA |
| Guide for C base editing | |
| rs112097551_C9 | GCAAGTAACGGATGCAGGGA |

**Table S7. Summary of PCR primers used in the present study**

| Gene symbol | Direction | Sequence (5'-3') |
|---|---|---|
| PCR primers for Sanger sequencing | | |
| RUVBL1 | Forward | ACTACTTACCAATGGCCCTG |
| | Reverse | GAGACAGAGAATCCCATGGG |
| RPN1 | Forward | GTAGGTCCTCAGAGCGCGTG |
| | Reverse | CAGAGTCATCCAAAATAAGG |
| rs112097551 | Forward | TCCTCTGTCCTTCCTTTCC |
| | Reverse | CATCTTGCCGATCTCTGAAC |
| Neutral locus | Forward | CCATGAGACAAGGAAGTAGTG |
| | Reverse | AGCAGTGGTGAGGAGAATA |
| Real-time qPCR primers | | |
| EEFSEC | Forward | GAGCGGCAAGTTCAAGAT |
| | Reverse | GTGGGTGTCGAAGACATAAC |
| GATA2 | Forward | TACAGCAGCGGACTCTT |
| | Reverse | GGTTCTGCCCATTCATCTT |
| RPN1 | Forward | ACCAGCCACCTCCTTATT |
| | Reverse | GGTCCACAAACCTCATCTTC |
| RAB7A | Forward | CCTAGATAGCTGGAGAGATGAG |
| | Reverse | CTGGTCTCAAAGTAGGGAATG |
| RUVBL1 | Forward | AAGGAGACCAAGGAAGTTTATG |
| | Reverse | CAGCTTCTACTCGCTCTTTC |
| GAPDH | Forward | ACCCAGAAGACTGTGGATGG |
| | Reverse | TTCAGCTCAGGGATGACCTT |

**Table S8. Lambda values in the single-variant association analyses in TOPMed**

| Trait | Lambda values | | |
|---|---|---|---|
| | **Unconditional analysis** | **Trait-specific conditional analysis** | **Trait-agnostic conditional analysis** |
| HCT | 1.021 | 1.020 | 1.015 |
| HGB | 1.019 | 1.019 | 1.015 |
| MCH | 1.036 | 1.034 | 1.029 |
| MCHC | 1.024 | 1.022 | 1.017 |
| MCV | 1.038 | 1.036 | 1.030 |
| RBC | 1.025 | 1.021 | 1.018 |
| RDW | 1.033 | 1.024 | 1.019 |

HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width.

**Table S9. Lead variants at the genome-wide significant loci of the marginal tests in TOPMed**
See Excel file.

**Table S10. Genome-wide significant variants at the 12 novel loci in the trait-specific conditional analysis in TOPMed**

See Excel file.

**Table S11. Ancestry-specific allele frequencies of the 14 novel lead variants at the 12 loci**

| Variant | Chr | Pos | Gene | Alternative allele | Reference allele | Alternative allele frequencies (%) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | European | African | Hispanic Latino | East Asian |
| rs112097551 | 3 | 1.3E+08 | *RNP1* | A | G | 0.069 | 0.940 | 0.400 | 0 |
| rs116635225 | 5 | 9.6E+07 | *ELL2* | A | G | 0.074 | 3.900 | 0.700 | 0 |
| rs986415672 | 10 | 1.3E+08 | *10q26* | T | C | 0.011 | 0 | 0 | 0 |
| rs11549407 | 11 | 5226774 | *HBB* | A | G | 0.016 | 0 | 0 | 0 |
| rs34598529 | 11 | 5227100 | *HBB* | C | T | 0 | 0.320 | 0 | 0 |
| rs535577177 | 11 | 7E+07 | *SHANK2* | A | G | 0 | 0 | 0.100 | 0 |
| rs370308370 | 14 | 1E+08 | *EIF5/MARK3* | A | G | 0 | 0 | 0 | 0.910 |
| rs868351380 | 16 | 55649 | *HBA1/HBA2* | C | G | 0.005 | 0 | 0.400 | 0 |
| rs372755452 | 16 | 199621 | *HBA1/HBA2* | A | AG | 0 | 0 | 0 | 1.100 |
| rs763477215 | 16 | 8.9E+07 | *PIEZO1* | A | ATCT | 0.355 | 0 | 0 | 0.050 |
| rs73494666 | 19 | 1253643 | *MIDN* | T | C | 0.614 | 51.7 | 4.700 | 0 |
| rs1368500441 | 19 | 2.9E+07 | *19q12* | A | G | 0.005 | 0 | 0 | 0 |
| rs228914 | 22 | 3.7E+07 | *TMPRSS6* | A | C | 88.7 | 96.6 | 79.2 | 99.8 |
| rs76723693 | X | 1.5E+08 | *G6PD* | G | A | 0 | 0.563 | 0.077 | 0 |

Chr, chromosome; Pos, position.

**Table S12. Replication results of the novel findings and the lead independent signals**
See Excel file.

**Table S13. Independent signals in the step-wise conditional analysis**
See Excel file.

**Table S14. Phenotypic variance explained by variants identified in the single variant association analysis**

| Trait | All | Known | Novel |
|---|---|---|---|
| HCT | 0.034 | 0.033 | 0.001 |
| HGB | 0.043 | 0.040 | 0.003 |
| MCH | 0.213 | 0.184 | 0.030 |
| MCHC | 0.047 | 0.041 | 0.006 |
| MCV | 0.179 | 0.153 | 0.028 |
| RBC | 0.126 | 0.117 | 0.010 |
| RDW | 0.118 | 0.109 | 0.009 |

HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCHC, mean corpuscular hemoglobin concentration; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width.

**Table S15. Summary of significant genes in the aggregated association analysis in TOPMed**
See Excel file.

**Table S16. Summary of significant genes in the aggregated association analysis adjusting for known and novel findings in TOPMed**

See Excel file.

**Table S17. Annotation of the rare variants identified in the aggregated analysis in TOPMed**
See Excel file.

**Table S18. Summary of pLoF and pKO variants in TOPMed freeze8 data** [1]

| Population | N [2] | No. of pLoF variants | No. of genes with at least one individual who is a pKO |
|---|---|---|---|
| African | 9,870 | 55,750 | 1,617 |
| Asian | 231 | 4,377 | 395 |
| European | 25,569 | 114,401 | 1,634 |
| Hispanic | 9,757 | 53,105 | 1,557 |

pLoF, predicted loss-of-function; pKO, predicted gene knockout; N, sample size.

1 No minor allele frequency filter was applied.

2 Sample sizes represented the number of individuals with blood-cell traits and genotype data available.

**Table S19. pLoF variants associated with RBC traits at P<1E-4 in TOPMed [1]**

| Population | Trait | Chr | Gene | rsID | Variant | MAF (%) | Type | Beta | SE | *P* |
|---|---|---|---|---|---|---|---|---|---|---|
| African | MCV | 2 | *WDSUB* | rs377262700 | chr2:159236041_G_A | 0.021 | stopgain | -2.682 | 0.592 | 6.09E-06 |
| African | RDW | 7 | *CD36* | rs3211938 | chr7:80671133_T_G[2] | 9.340 | stopgain | 0.244 | 0.043 | **1.24E-08** |
| Hispanic | MCV | 4 | *SNX25* | rs1200775460 | chr4:185339389_AG_A | 0.022 | frameshift | 2.441 | 0.543 | 7.08E-06 |
| Hispanic | MCH | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.022 | stopgain | -2.611 | 0.505 | **2.36E-07** |
| Hispanic | MCV | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.022 | stopgain | -2.970 | 0.506 | **4.58E-09** |
| Hispanic | RBC | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.022 | stopgain | 2.272 | 0.506 | 7.18E-06 |
| European | HCT | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.018 | stopgain | -1.545 | 0.333 | 3.39E-06 |
| European | HGB | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.018 | stopgain | -2.052 | 0.332 | **6.65E-10** |
| European | MCH | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.017 | stopgain | -2.974 | 0.494 | **1.73E-09** |
| European | MCV | 11 | *HBB* | rs11549407 | chr11:5226774_G_A[3] | 0.015 | stopgain | -2.981 | 0.494 | **1.66E-09** |
| European | HCT | 11 | *CD6* | rs759187282 | chr11:61017803_G_T | 0.006 | stopgain | -2.573 | 0.575 | 7.73E-06 |
| European | HGB | 11 | *CD6* | rs759187282 | chr11:61017803_G_T | 0.006 | stopgain | -2.575 | 0.574 | 7.39E-06 |
| Meta-analysis | RDW | 1 | *SMIM1* | rs566629828 | chr1:3775433_AGTCAGCCTAGGGGCTGT_A[4] | 1.610 | frameshift | 0.303 | 0.068 | 8.22E-06 |
| Meta-analysis | MCV | 18 | *SERPINB11* | rs760239610 | chr18:63712688_C_CATCAGGTA | 0.150 | frameshift | -0.681 | 0.150 | 5.60E-06 |

pLoF, predicted loss-of-function; HCT, hematocrit; HGB, hemoglobin; MCH, mean corpuscular hemoglobin; MCV, mean corpuscular volume; RBC, red blood cell count; RDW, red blood cell width; Chr, chromosome; MAF, minor allele frequency.

1 The P values of pLoF variants that reached genome-wide significance were in bold (African: P<8.97E-7; Hispanic: P<9.42E-7; European: P<4.37E-7).

2 Well known CD36 null allele.

3 Well known beta-thalassemia allele.

4 This frameshift indel is responsible for the Vel blood group.

**Table S20. pKO variants associated with RBC traits at P<1E-4 in TOPMed [1]**

| Population | Trait | Chr | Gene | rsID | Variants | MAF (%) | Type | N Total | N KO | Beta | SE | *P* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| African | MCH | 7 | *ZNF3* | rs777843966 | chr7:100064797_GT_G | 0.008 | frameshift | 6042 | 200 | 0.277 | 0.070 | 8.53E-05 |
| | | | | rs987730433 | chr7:100064875_C_CA | 0.008 | frameshift | | | | | |
| | | | | rs71689664 | chr7:100064888_GTAGT_G | 18.3 | frameshift | | | | | |
| | | | | rs745468385 | chr7:100071151_ACT_A | 0.008 | frameshift | | | | | |
| | | | | rs774923137 | chr7:100071181_TG_T | 0.008 | frameshift | | | | | |
| | | | | rs988854061 | chr7:100079535_C_T | 0.008 | splicing | | | | | |
| African | MCV | 7 | *ZNF3* | rs777843966 | chr7:100064797_GT_G | 0.007 | frameshift | 7198 | 239 | 0.267 | 0.064 | 3.60E-05 |
| | | | | rs987730433 | chr7:100064875_C_CA | 0.007 | frameshift | | | | | |
| | | | | rs71689664 | chr7:100064888_GTAGT_G | 18.4 | frameshift | | | | | |
| | | | | rs745468385 | chr7:100071151_ACT_A | 0.007 | frameshift | | | | | |
| | | | | rs774923137 | chr7:100071181_TG_T | 0.007 | frameshift | | | | | |
| | | | | rs988854061 | chr7:100079535_C_T | 0.007 | splicing | | | | | |

pKO, predicted gene knockout; MCH, mean corpuscular hemoglobin; MCV, mean corpuscular volume; Chr, chromosome; MAF, minor allele frequency.

1 No pKO variant reached genome-wide significance (African: P<3.09E-5; Hispanic: P<3.21E-5; European: P<3.06E-5).

**Supplemental Methods**

**Participating studies**

*Amish*

The Amish Complex Disease Research Program includes a set of large community-based studies focused largely on cardiometabolic health carried out in the Old Order Amish (OOA) community of Lancaster, Pennsylvania (http://medschool.umaryland.edu/endocrinology/amish/research-program.asp). The OOA population of Lancaster County, PA immigrated to the Colonies from Western Europe in the early 1700's. There are now over 30,000 OOA individuals in the Lancaster area, nearly all of whom can trace their ancestry back 12-14 generations to approximately 700 founders. Investigators at the University of Maryland School of Medicine have been studying the genetic determinants of cardiometabolic health in this population since 1993. To date, over 7,000 Amish adults have participated in one or more of our studies.

Due to their ancestral history, the OOA are enriched for rare exonic variants that arose in the population from a single founder (or small number of founders) and propagated through genetic drift. Many of these variants have large effect sizes and identifying them can lead to new biological insights about health and disease. The parent study for this WGS project provides one (of multiple) examples. In our parent study, we identified through a genome-wide association analysis a haplotype that was highly enriched in the OOA that is associated with very high LDL-cholesterol levels. At the present time, the identity of the causative SNP – and even the implicated gene – is not known because the associated haplotype contains numerous genes, none of which are obvious lipid candidate genes. A major goal of the WGS that will be obtained through the NHLBI TOPMed Consortium will be to identify functional variants that underlie some of the large effect associations observed in this unique population.

*ARIC*

The ARIC study is a population-based cohort study consisting of 15,792 men and women that were drawn from four U.S. communities (Suburban Minneapolis, Minnesota; Washington County, Maryland; Forsyth County, North Carolina, and Jackson, Mississippi)[1]. It was designed to investigate the causes of atherosclerosis and its clinical outcomes, and variation in cardiovascular risk factors, medical care, and disease by race, sex, location, and date. For TOPMed WGS, the study over-sampled participants with incident VTE. Participants were between age 45 and 64 years at their baseline examination in 1987-1989 when blood was drawn for DNA extraction and participants consented to genetic testing.

*BioMe*

The Charles Bronfman Institute for Personalized Medicine at Mount Sinai Medical Center (MSMC), BioMe Biobank, founded in September 2007, is an ongoing, broadly-consented electronic health record-linked

clinical care biobank that enrolls participants non-selectively from the Mount Sinai Medical Center patient population. The MSMC serves diverse local communities of upper Manhattan, including Central Harlem (86% African American), East Harlem (88% Hispanic/Latino), and Upper East Side (88% Caucasian/White) with broad health disparities.

*CARDIA*

The Coronary Artery Risk Development in Young Adults (CARDIA) Study is a study examining the development and determinants of clinical and subclinical cardiovascular disease and their risk factors. It began in 1985-1986 with a group of 5,115 black and white men and women aged 18-30 years. The participants were selected so that there would be approximately the same number of people in subgroups of race, gender, education (high school or less and more than high school) and age (18-24 and 25-30) in each of 4 centers: Birmingham, AL; Chicago, IL; Minneapolis, MN; and Oakland, CA.

*CHS*

The Cardiovascular Health Study (CHS) is a population-based cohort study of risk factors for coronary heart disease and stroke in adults 65 years and older conducted across four field centers [2]. The original predominantly European ancestry cohort of 5,201 persons was recruited in 1989-1990 from random samples of people on Medicare eligibility lists from four US communities. Subsequently, an additional predominantly African-American cohort of 687 persons was enrolled for a total sample of 5,888. Institutional review committees at each field center approved the CHS, and participants gave informed consent. Blood samples were drawn from all participants at their baseline examination, and DNA was subsequently extracted from available samples. These analyses were limited to participants with available DNA who also consented to genetic studies. Participants were examined annually from enrollment to 1999 and continued to be under surveillance for stroke following 1999.

*COPDGene*

COPDGene (also known as the Genetic Epidemiology of COPD Study) is an NIH-funded, multicenter study. A study population of more than 10,000 smokers (1/3 African American and 2/3 non-Hispanic White) has been characterized with a study protocol including pulmonary function tests, chest CT scans, six minute walk testing, and multiple questionnaires. Five years after this initial visit, all available study participants are being brought back for a follow-up visit with a similar study protocol. This study has been used for epidemiologic and genetic studies. Previous genetic analysis in this study has been based on genome-wide SNP genotyping data. Approximately 1,900 subjects underwent whole genome sequencing

in this NHLBI WGS project, including severe COPD subjects and non-COPD smoking controls. The COPDGene Study web site is: http://www.copdgene.org/.

*FHS*

FHS is a three-generation, single-site, community-based, ongoing cohort study that was initiated in 1948 to investigate prospectively the risk factors for CVD including stroke. It now comprises 3 generations of participants: the Original cohort followed since 1948 [3]; their Offspring and spouses of the Offspring, followed since 1971 [4]; and children from the largest Offspring families enrolled in 2002 (Gen 3) [5]. The Original cohort enrolled 5,209 men and women who comprised two-thirds of the adult population then residing in Framingham, MA. Survivors continue to receive biennial examinations. The Offspring cohort comprises 5,124 persons (including 3,514 biological offspring) who have been examined approximately once every 4 years. The Gen 3 cohort contains 4,095 participants.

*GeneSTAR*

In 1982 The Johns Hopkins Sibling and Family Heart Study was created to study patterns of coronary heart disease and related risk factors in families with early-onset coronary disease, identified from 10 Baltimore area Hospitals. GeneSTAR continues to study mechanisms of coronary heart disease and stroke in families using novel models and exciting new methods. GeneSTAR is a family-based study in initially healthy brothers and sisters, and offspring of people with early-onset coronary disease, The goal is to discover and amplify mechanisms of stroke and coronary heart disease. Our African American and European American family cohort has undergone extensive screening, genetic testing, and follow-up for new cardiovascular disease, stroke, and other clinical events for 5 to 32 years.

*HCHS/SOL*

The Hispanic Community Health Study/Study of Latinos (HCHS/SOL) is a multi-center study of Hispanic/Latino populations with the goal of determining the role of acculturation in the prevalence and development of diseases, and to identify other traits that impact Hispanic/Latino health [6] . The study is sponsored by the National Heart, Lung, and Blood Institute (NHLBI) and other institutes, centers, and offices of the National Institutes of Health (NIH). Recruitment began in 2006 with a target population of 16,000 persons of Cuban, Puerto Rican, Dominican, Mexican or Central/South American origin. Participants were recruited through four sites affiliated with San Diego State University, Northwestern University in Chicago, Albert Einstein College of Medicine in Bronx, New York, and the University of Miami. Recruitment was implemented through a two-stage area household probability design [6]. The study enrolled 16,415 participants who were self-identified Hispanic/Latino and aged 18-74 years and the

extensive psycho-social and clinical assessments were conducted during 2008-2011. Annual telephone follow-up interviews are ongoing since study inception. During the 2014-2017 second visit, the participants were re-examined again of various health outcomes of interest.

*JHS*

The Jackson Heart Study (JHS, https://www.jacksonheartstudy.org/jhsinfo/) is a large, community-based, observational study whose participants were recruited from urban and rural areas of the three counties (Hinds, Madison and Rankin) that make up the Jackson, MS metropolitan statistical area (MSA). Participants were enrolled from each of 4 recruitment pools: random, 17%; volunteer, 30%; currently enrolled in the Atherosclerosis Risk in Communities (ARIC) Study, 31% and secondary family members, 22%. Recruitment was limited to non-institutionalized adult African Americans 35-84 years old, except in a nested family cohort where those 21 to 34 years of age were also eligible. The final cohort of 5,301 participants included 6.59% of all African American Jackson MSA residents aged 35-84 during the baseline exam (N-76,426, US Census 2000). Among these, approximately 3,700 gave consent that allows genetic research and deposition of data into dbGaP. Major components of three clinic examinations (Exam 1 – 2000-2004; Exam 2 – 2005-2008; Exam 3 – 2009-2013) include medical history, physical examination, blood/urine analytes and interview questions on areas such as: physical activity; stress, coping and spirituality; racism and discrimination; socioeconomic position; and access to health care. Extensive clinical phenotyping includes anthropometrics, electrocardiography, carotid ultrasound, ankle-brachial blood pressure index, echocardiography, CT chest and abdomen for coronary and aortic calcification, liver fat, and subcutaneous and visceral fat measurement, and cardiac MRI. At 12-month intervals after the baseline clinic visit (Exam 1), participants have been contacted by telephone to: update information; confirm vital statistics; document interim medical events, hospitalizations, and functional status; and obtain additional sociocultural information. Questions about medical events, symptoms of cardiovascular disease and functional status are repeated annually. Ongoing cohort surveillance includes abstraction of medical records and death certificates for relevant International Classification of Diseases (ICD) codes and adjudication of nonfatal events and deaths. CMS data are currently being incorporated into the dataset.

*MESA*

The MESA study is a study of the characteristics of subclinical cardiovascular disease (disease detected non-invasively before it has produced clinical signs and symptoms) and the risk factors that predict progression to clinically overt cardiovascular disease or progression of the subclinical disease [7]. MESA researchers study a diverse, population-based sample of 6,814 asymptomatic men and women aged 45-84. Thirty-eight percent of the recruited participants are white, 28 percent African-American, 22 percent

Hispanic, and 12 percent Asian, predominantly of Chinese descent. Participants were recruited from six field centers across the United States: Wake Forest University, Columbia University, Johns Hopkins University, University of Minnesota, Northwestern University and the University of California - Los Angeles.

*SAFS*

The San Antonio Family Study (SAFS) is a complex pedigree-based mixed longitudinal study designed to identify low frequency or rare variants influencing susceptibility to cardiovascular disease, using WGS information from 2,590 individuals in large Mexican American pedigrees from San Antonio, Texas. The major objectives of this study are to identify low frequency or rare variants in and around known common variant signals for CVD, as well as to find novel low frequency or rare variants influencing susceptibility to CVD.

*WHI*

The Women's Health Initiative (WHI) is a long-term, prospective, multi-center cohort study that investigates post-menopausal women's health [8]. WHI was funded by the National Institutes of Health and the National Heart, Lung, and Blood Institute to study strategies to prevent heart disease, breast cancer, colon cancer, and osteoporotic fractures in women 50-79 years of age. WHI involves 161,808 women recruited between 1993 and 1998 at 40 centers across the US. The study consists of two parts: the WHI Clinical Trial which was a randomized clinical trial of hormone therapy, dietary modification, and calcium/Vitamin D supplementation, and the WHI Observational Study, which focused on many of the inequities in women's health research and provided practical information about the incidence, risk factors, and interventions related to heart disease, cancer, and osteoporotic fractures. For TOPMed WGS, the study over-sampled participants with incident stroke and VTE. The remaining samples were age- and ethnicity-matched controls without stroke or VTE.

**Phenotype harmonization**

Because multiple studies contributed to the analysis, RBC phenotypes were harmonized across studies such that they could be analyzed together (https://www.biorxiv.org/content/10.1101/2020.06.18.146423v1). For most studies, variables were obtained from dbGaP. Data for the BioME study, the COPDGene study, and some participants in the WHI study were transferred directly to investigators. The study variables were QC'd to identify recording errors or other problematic data. After QC, variables were converted to a consistent measurement unit across studies. When possible, harmonized variables were calculated using study-derived variables. If QC uncovered data quality issues with the study-derived variable or if the study

did not provide the specific variable of interest, the harmonized variable was instead calculated from the components (e.g., MCH = 10 * hemoglobin / red cell count). Finally, QC of the harmonized variable was performed to check that no large differences between studies remained.

Both FHS and ARIC measured phenotypes at multiple times in a given participant. For these studies, a single measurement for each participant was selected for each trait. To maximize the sample size and minimize batch effects within studies, harmonized data for different RBC traits for a single subject may have been measured at different visits. Within FHS, only the Offspring cohort had measurements at multiple exams. For these participants, the measurement at the most recent exam was chosen. For the ARIC study, the visit with the most non-missing phenotype values across all participants was chosen first. For subjects without measurements at this visit, the visit with the next most non-missing values was chosen, and so forth. As a consequence, values for the same participant for different RBC phenotypes were sometimes measured at different visits."

Trait-specific QC procedures were also performed. We excluded participants with HCT values >80% and those with HCT <5% (n=14). Similarly for HGB, we excluded participants with HGB measurements >30 g/dL and <5 g/dL (n=31). For participants from WHI, for both the HCT and HGB analyses, we excluded those with an HCT/HGB ratio >7 (n=11). Participants with MCH values >75 pg were excluded from analysis (n=2). In MCHC, we excluded participants with measurements ≥60 g/dL (n=1). For the MCV analysis, outliers with values >150 fL were excluded (n=2). In the RBC and RDW analyses, no participants were excluded based upon measured trait values.

### Statistical analyses

*Genetic Ancestry and Relatedness*

Principal components (PCs) of genetic ancestry and pairwise relatedness measures were estimated for all 140,062 samples included in the TOPMed 'freeze 8' genotype release. Autosomal genetic variants passing the quality filter with a MAF > 0.01 and missing call rate < 0.01 were LD-pruned with an $r^2$ threshold of 0.1 to obtain a set of 638,486 effectively independent variants for genetic ancestry and relatedness estimation. PC-AiR [9] was used to obtain ancestry informative PCs robust to familial relatedness; the first 11 PCs showed evidence of population structure. PC-Relate [10] was then used to estimate pairwise kinship coefficients (KCs) for all pairs of samples, conditional on the genetic ancestry captured by PC-AiR PCs 1-11; these KC estimates reflect only recent genetic relatedness, e.g. due to pedigree structure. The PC-Relate KC estimates were used to construct a 4th degree sparse, block-diagonal, empirical kinship matrix (KM) for association testing, using the procedure recommended in Gogarten et al [11].: any pair of samples with estimated KC > $2^{(-11/2)}$ ~ 0.022 were clustered in the same block; all KC estimates within a block of samples were kept, regardless of value; and all KC estimates between blocks were set to 0. By using a sparse block-

diagonal KM, the association tests are more computationally efficient yet recent genetic relatedness is still accounted for. We subset the freeze-wide PCs and sparse KM to the appropriate set of participants for each analysis.

*HARE for Imputation of Race/Population Membership using Genetic Ancestry*

Ancestry groups were based on a combination of participants reported race/ethnicity and genetic ancestry represented by PCs from PC-AiR[9]. To infer race/population group membership for participants with missing values, we used the HARE method [12] . HARE is a machine learning algorithm that uses a support vector machine (SVM) to determine stratum assignment, taking as input genetically estimated PC values and reported race/ethnicity for each participant. Strata are defined by the unique reported race/ethnicity values provided, then the HARE SVM uses the input (training) data to learn the probability of stratum membership across the entire PC space. The output of HARE consists of multinomial probability vectors of stratum membership for each participant. HARE was run on a subset of samples included in the TOPMed freeze 8 genotype release; specifically, samples for participants from non-US populations (e.g. Costa Rica) and the Amish participants (because they were very distinct in PC space) were excluded from the HARE analysis. HARE was run using the first 9 PC-AiR PCs generated on this subset of samples to represent genetic ancestry with the following reported race/population groups: Asian, Black, Central American, Cuban, Dominican, Mexican, Puerto Rican, South American, and White. The genetic data from the 31,918 participants with either unreported or non-specific (e.g. 'Multiple' or 'Other') race and population membership was included in the HARE analysis, but they were not used to train the SVM. These participants were assigned to a population stratum based on their highest HARE output probability of membership. All other participants remained in the population stratum corresponding to their reported race/population group. Amish participants were assigned to their own stratum.

*Fitting the Linear Mixed Model*

The linear mixed model (LMM) can be written as $Y = G\beta + X\alpha + \epsilon$, where $Y$ is the $(n \ x \ 1)$ vector of outcome values; $G$ is an $(n \ x \ m)$ matrix of alternate allele counts for each of the $n$ individuals at the $m$ variants of interest ($m = 1$ for a single variant analysis) with effect sizes given by the $(m \ x \ 1)$ vector $\beta$; $X$ is the $(n \ x \ k)$ matrix of fixed effect covariates including an intercept with effect sizes given by the $(k \ x \ 1)$ vector $\alpha$; and $\epsilon \sim N(0, \Sigma)$ is the $(n \ x \ 1)$ vector of errors with covariance matrix $\Sigma$ that captures both genetic covariance due to relatedness/kinship and residual variance structure. Given the true $\Sigma$, we could estimate $\beta$ using generalized least squares (GLS). However, we can simplify this GLS problem to an ordinary least squares (OLS) problem by pre-multiplying both sides of the equation by the matrix $C$, the Cholesky-decomposition of $\Sigma^{-1}$, such that $C'C = \Sigma^{-1}$ and $C'\Sigma C = I$, where $I$ is the $(n \ x \ n)$ identity matrix. Further,

by the Frisch-Waugh Lovell theorem[13], we can adjust for the covariates in the new OLS model, $CX$, by pre-multiplying $CY$ and $CG$ by the annihilator matrix $[I - (CX)((CX)'(CX))^{-1}(CX)']$. Ultimately, the original GLS problem can be re-written as the linear regression model $Y^* = G^*\beta + \epsilon^*$, where $Y^* = MY$, $G^* = MG$, and $M = [I - CX(X'C'CX)^{-1}X'C']C$. In practice, we use REML to estimate $\hat{\Sigma}$ under the null hypothesis that $\beta = 0$ (i.e. fit the null model) and calculate the estimate of the matrix $\hat{M}$.

*Score Tests and Approximate Variant Effect Sizes*

Given $Y^*$ and $G^*$, a joint score test for the set of $m$ variants can be performed, where the score is $U = G^{*'}Y^*$, the variance of the score is $V = G^{*'}G^*$, and the test statistic is $T_G = U'V^{-1}U \sim \chi^2_m$. The score and the Wald tests are approximately asymptotically equivalent when $\beta$ is small (as is typical for GWAS), so the variant effect sizes can be reasonably approximated from the score test as $\hat{\beta} \approx V^{-1}U = (G^{*'}G^*)^{-1}(G^{*'}Y^*)$, and their covariance matrix can be reasonably approximated from the score test as $\widehat{var}(\hat{\beta}) \approx V^{-1} = (G^{*'}G^*)^{-1}$.[14] Note that these are the score tests used for the single variant association analysis, where each variant genome-wide is tested individually (i.e. $m = 1$).

*Proportion of Variance Explained Jointly by a set of variants*

To estimate the proportion of phenotypic variance explained (PVE) by the $m$ variants in $G$, we use the formula $PVE = 1 - RSS_1/RSS_0$, where $RSS_0$ and $RSS_1$ are the residual sums of squares computed from the null model, and the model including the $m$ variants of interest, respectively. Under the null model, we have that $RSS_0 = Y^{*'}Y^*$, and from the model $Y^* = G^*\beta + \epsilon^*$, we have that $RSS_1 = (Y^* - G^*\hat{\beta})'(Y^* - G^*\hat{\beta})$. Using the approximation for $\hat{\beta}$ given above, we get that $RSS_1 \approx Y^{*'}Y^* - Y^{*'}G^*(G^{*'}G^*)^{-1}G^{*'}Y^* = Y^{*'}Y^* - T_G$, and the estimate $\widehat{PVE} \approx T_G/(Y^{*'}Y^*)$. It's worth noting that using this approach to estimate the PVE for the set of $m$ variants jointly should provide a more accurate estimate than estimating the PVE for each variant separately and summing, as this joint approach accounts for the covariance between the variant effect sizes, as measured by $V^{-1}$ (the separate approach is equivalent to $\widehat{PVE} = [U'diag(V)^{-1}U]/(Y^{*'}Y^*)$, where $diag(V)$ is an $(m \, x \, m)$ matrix of just the diagonal of the $V$ matrix). This joint PVE calculation is implemented in the GENESIS software [11] with the jointScoreTest function.

*Conditional analyses*

We performed three types of conditional analysis in the discovery stage. The first conditional analyses adjusted each trait for variants that were previously reported to be associated with the particular RBC trait and that passed the QC filter. These variants were pruned to a set with linkage disequilibrium (LD) r2 < 0.8 such that a variant with a more significant p-value was preferentially retained over those with higher p-

values. Known variants failing the QC filter were included if any variants within 1 MB of the known variant remained significant after adjusting for the passing variants only. We refer to this analysis as the "RBC trait-specific conditional analysis". In the second conditional analysis, we included all previously reported variants for *any* of the seven RBC traits as well as any failed variants included in any of the trait-specific conditional analyses. Variants were again pruned to LD $r^2 < 0.8$ with preferential selection based on p-value. We refer to the second conditional analysis as the "RBC trait-agnostic conditional analysis". Finally, we performed iterative conditional analysis by chromosome for each trait to identify an independent set of associated variants. For this third conditional analysis, we started with the association results from the trait-specific conditional analysis. For each chromosome, we identified the most significant variant (if any, using a $5\times10^{-9}$ threshold) as the 'peak variant' and then fit a new null model adjusted for both the previous set of conditional variants from the trait-specific conditional analyses as well as this peak variant, and calculated new score test statistics. If any variant was significant at the $5\times10^{-9}$ level in the new score tests (regardless of its significance level in the original trait-specific conditional results), we performed a second round of conditional analysis, re-estimating the null model and calculating the score test statistics, adjusting for the new peak variant along with the original trait-specific conditional variants and the first peak variant. We continued this procedure iteratively, adding any new 'peak variants' into the list of variants to condition on, re-fitting the null model, and calculating the updated score statistics, until no additional variants were significant at the $5\times10^{-9}$ level. Finally, the variants identified across all chromosomes in this iterative conditional analysis were combined into a set of "conditionally-independent variants" for each trait.

*Aggregation Strategies*

For aggregate association testing, five distinct methods were used to aggregate rare variants into gene-based groups using GENCODE v29 gene model. Three strategies only included coding variants, and two strategies additionally included non-coding variants. Variants were further filtered using one or more deleterious prediction scores to enrich for likely causal variants. The detail method used for each strategy is provided below

1. **Coding filter 1 - Stringent (C1-S)**: This strategy includes high confidence predicted LoF variants inferred using LOFTEE (https://github.com/konradjk/loftee), missense variants predicted deleterious by all of SIFT4G[26633127]<=0.05, Polyphen2_HDIV>0.5[20354512], Polyphen2_HVAR>0.5[20354512], and variants predicted as "Deleterious" by LRT [19602639] and inframe indels or synonymous variants with Fathmm-XF score[28968714] > 0.5

2. **Coding filter 1 - Relaxed (C1-R)**: This strategy is same as C1-S but the missense filter was relaxed to retain variants predicted deleterious by any of SIFT4G, Polyphen2_HDIV, Polyphen2_HVAR scores

3. **Coding filter 2 - Relaxed (C2-R)**: This strategy is the same as C1-S but missense variants were filtered using MetaSVM score and a relatively relaxed set of missense variants was retained by applying MetaSVM score [25552646] > 0 filter

4. **Coding filter 2 - Relaxed & Non-coding filter-Relaxed (C2-R+NC-R)**: This strategy includes variants included in C2-R and additional regulatory variants. Regulatory variants were included if they overlapped with enhancer(s) or promoters linked to a gene using GeneHancer[28605766], or 5 Kb upstream of the Transcription start site. Within these regions only those variants were retained which had Fathmm-XF score > 0.5 **or** overlap with regions labelled as either"CTCF binding sites," "Transcription factor binding sites" as annotated by the Ensembl regulatory build annotation[25887522]

5. **Coding filter 2 - Relaxed & Non-coding filter-Stringent (C2-R +NC-S)** :This strategy includes variants included in C2-R and additional regulatory variants using a stringent filtering criteria. Even in this method regulatory variants in Genehancer[28605766] linked regulatory regions and 5 Kb upstream of the Transcription start site of gene were included. However within these regions only those variants were retained which had Fathmm-XF score > 0.5 **and** which overlapped with regions labelled as "Promoters," "Promoter flanking regions," "Enhancers," "CTCF binding sites," "Transcription factor binding sites" or "Open chromatin regions" as annotated by the Ensembl regulatory build annotation[25887522].

The annotation based variant filtering and gene based aggregation was performed using TOPMed freeze 8 WGSA Google BigQuery annotation database on the BiodataCatalyst powered by Seven Bridges platform (http://doi.org/10.5281/zenodo.3822858). The annotation database was built using variant annotations generated by Whole genome Sequence annotator version v0.8 [26395054] and formatted by WGSAParsr version 6.3.8 (https://github.com/UW-GAC/wgsaparsr). The GENCODE v29 gene model based varint consequences were obtained from Ensembl Variant effect predictor (VEP)[26683364] incorporated within WGSA. When using a deleteriousness prediction score, respective author recommended cut points were used to retain likely deleterious variants.

**References**

1. (1989). The atherosclerosis risk in communities (aric) study: design and objectives. Am. J. Epidemiol. *129*, 687–702.

2. Fried, L.P., Borhani, N.O., Enright, P., Furberg, C.D., Gardin, J.M., Kronmal, R.A., Kuller, L.H., Manolio, T.A., Mittelmark, M.B., and Newman, A. (1991). The Cardiovascular Health Study: design and rationale. Ann. Epidemiol. *1*, 263–276.

3. Dawber, T.R., and Kannel, W.B. (1966). The Framingham study. An epidemiological approach to coronary heart disease. Circulation *34*, 553–555.

4. Feinleib, M., Kannel, W.B., Garrison, R.J., McNamara, P.M., and Castelli, W.P. (1975). The Framingham Offspring Study. Design and preliminary data. Prev. Med. *4*, 518–525.

5. Splansky, G.L., Corey, D., Yang, Q., Atwood, L.D., Cupples, L.A., Benjamin, E.J., D'Agostino, R.B., Fox, C.S., Larson, M.G., Murabito, J.M., et al. (2007). The Third Generation Cohort of the National Heart, Lung, and Blood Institute's Framingham Heart Study: design, recruitment, and initial examination. Am. J. Epidemiol. *165*, 1328–1335.

6. Lavange, L.M., Kalsbeek, W.D., Sorlie, P.D., Avilés-Santa, L.M., Kaplan, R.C., Barnhart, J., Liu, K., Giachello, A., Lee, D.J., Ryan, J., et al. (2010). Sample design and cohort selection in the Hispanic Community Health Study/Study of Latinos. Ann. Epidemiol. *20*, 642–649.

7. Bild, D.E., Bluemke, D.A., Burke, G.L., Detrano, R., Diez Roux, A.V., Folsom, A.R., Greenland, P., Jacob, D.R., Kronmal, R., Liu, K., et al. (2002). Multi-Ethnic Study of Atherosclerosis: objectives and design. Am. J. Epidemiol. *156*, 871–881.

8. (1998). Design of the Women's Health Initiative clinical trial and observational study. The Women's Health Initiative Study Group. Control. Clin. Trials *19*, 61–109.

9. Conomos, M.P., Miller, M.B., and Thornton, T.A. (2015). Robust inference of population structure for ancestry prediction and correction of stratification in the presence of relatedness. Genet. Epidemiol. *39*, 276–293.

10. Conomos, M.P., Reiner, A.P., Weir, B.S., and Thornton, T.A. (2016). Model-free Estimation of Recent Genetic Relatedness. Am. J. Hum. Genet. *98*, 127–148.

11. Gogarten, S.M., Sofer, T., Chen, H., Yu, C., Brody, J.A., Thornton, T.A., Rice, K.M., and Conomos, M.P. (2019). Genetic association testing using the GENESIS R/Bioconductor package. Bioinformatics *35*, 5346–5348.

12. Fang, H., Hui, Q., Lynch, J., Honerlaw, J., Assimes, T.L., Huang, J., Vujkovic, M., Damrauer, S.M., Pyarajan, S., Gaziano, J.M., et al. (2019). Harmonizing Genetic Ancestry and Self-identified Race/Ethnicity in Genome-wide Association Studies. Am. J. Hum. Genet. *105*, 763–772.

13. Frisch, R., & Waugh, F. V. (1933). Partial time regressions as compared with individual trends.

Econometrica: Journal of the Econometric Society, 387-401.

14. Zhou, B., Shi, J., and Whittemore, A.S. (2011). Optimal methods for meta-analysis of genome-wide association studies. Genet. Epidemiol. *35*, 581–591.

## Acknowledgements

| TOPMed Accession # | TOPMed Project | Parent Study Name | TOPMed Phase | Omics Center | Omics Support |
|---|---|---|---|---|---|
| phs000956 | Amish | Amish | 1 | Broad Genomics | 3R01HL121007-01S1 |
| phs001211 | AFGen | ARIC AFGen | 1 | Broad Genomics | 3R01HL092577-06S1 |
| phs001211 | VTE | ARIC | 2 | Baylor | 3U54HG003273-12S2 / HHSN268201500015C |
| phs001644 | AFGen | BioMe AFGen | 2.4 | MGI | 3UM1HG008853-01S2 |
| phs001644 | BioMe | BioMe | 3 | Baylor | HHSN268201600033I |
| phs001644 | BioMe | BioMe | 3 | MGI | HHSN268201600037I |
| phs001612 | CARDIA | CARDIA | 3 | Baylor | HHSN268201600033I |
| phs001368 | CHS | CHS | 3 | Baylor | HHSN268201600033I |
| phs001368 | VTE | CHS VTE | 2 | Baylor | 3U54HG003273-12S2 / HHSN268201500015C |
| phs000951 | COPD | COPDGene | 1 | NWGC | 3R01HL089856-08S1 |
| phs000951 | COPD | COPDGene | 2 | Broad Genomics | HHSN268201500014C |
| phs000951 | COPD | COPDGene | 2.5 | Broad Genomics | HHSN268201500014C |
| phs000974 | AFGen | FHS AFGen | 1 | Broad Genomics | 3R01HL092577-06S1 |
| phs000974 | FHS | FHS | 1 | Broad Genomics | 3U54HG003067-12S2 |
| phs001218 | AA_CAC | GeneSTAR AA_CAC | 2 | Broad Genomics | HHSN268201500014C |

| phs001218 | GeneSTAR | GeneSTAR | legacy | Illumina | R01HL112064 |
|-----------|----------|----------|--------|----------|-------------|
| phs001218 | GeneSTAR | GeneSTAR | 2 | Psomagen | 3R01HL112064-04S1 |
| phs001395 | HCHS_SOL | HCHS_SOL | 3 | Baylor | HHSN268201600033I |
| phs000964 | JHS | JHS | 1 | NWGC | HHSN268201100037C |
| phs001416 | AA_CAC | MESA AA_CAC | 2 | Broad Genomics | HHSN268201500014C |
| phs001416 | MESA | MESA | 2 | Broad Genomics | 3U54HG003067-13S1 |
| phs001215 | SAFS | SAFS | 1 | Illumina | 3R01HL113323-03S1 |
| phs001215 | SAFS | SAFS | legacy | Illumina | R01HL113322 |
| phs001237 | WHI | WHI | 2 | Broad Genomics | HHSN268201500014C |