





# Generalization and fine mapping of red blood cell trait genetic associations to multi-ethnic populations: The PAGE study

Chani J. Hodonsky<sup>1</sup>  | Claudia Schurmann<sup>2,3</sup> | Ursula M. Schick<sup>2,3,4</sup> | Jonathan Kocarnik<sup>4</sup>  | Ran Tao<sup>5</sup> | Frank J. A. van Rooij<sup>6</sup> | Christina Wassel<sup>7</sup> | Steve Buyske<sup>8</sup> | Myriam Fornage<sup>9</sup> | Lucia A. Hindorff<sup>10</sup> | James S. Floyd<sup>11,12</sup> | Santhi K. Ganesh<sup>13,14</sup> | Dan-Yu Lin<sup>15</sup> | Kari E. North<sup>1</sup> | Alex P. Reiner<sup>4,12</sup> | Ruth J. F. Loos<sup>2,3</sup> | Charles Kooperberg<sup>4</sup> | Christy L. Avery<sup>1,16</sup>

<sup>1</sup>Department of Epidemiology, University of North Carolina Gillings School of Public Health, Chapel Hill, North Carolina; <sup>2</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York; <sup>3</sup>The Genetics of Obesity and Related Metabolic Traits Program, The Icahn School of Medicine at Mount Sinai, New York; <sup>4</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, Washington; <sup>5</sup>Department of Biostatistics, Vanderbilt University Medical Center, Nashville, Tennessee; <sup>6</sup>Department of Epidemiology, Erasmus University Medical Center, Rotterdam, 3000, the Netherlands; <sup>7</sup>Department of Pathology and Laboratory Medicine, College of Medicine, University of Vermont, Burlington, Vermont; <sup>8</sup>Department of Statistics and Biostatistics, Hill Center, Rutgers, The State University of New Jersey, 110 Frelinghuysen Rd., Piscataway, New York; <sup>9</sup>Institute of Molecular Medicine and Human Genetics Center, University of Texas Health Science Center at Houston, Houston, Texas; <sup>10</sup>Division of Genomic Medicine, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland; <sup>11</sup>Departments of Medicine, University of Washington, Seattle, Washington; <sup>12</sup>Department of Epidemiology, University of Washington, Seattle, Washington; <sup>13</sup>Department of Human Genetics, University of Michigan Medical School, Ann Arbor, Michigan; <sup>14</sup>Department of Internal Medicine, University of Michigan Medical School, Ann Arbor, Michigan; <sup>15</sup>Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina; <sup>16</sup>Carolina Population Center, University of North Carolina, Chapel Hill, North Carolina

## Correspondence

Chani J. Hodonsky, 123 W Franklin St, St 4208-A, Chapel Hill, NC 27514.  
Email: chani\_hodonsky@unc.edu

## Funding information

National Heart, Lung, and Blood Institute; NIH; and U.S. Department of Health and Human Services, Grant numbers: N01WH22110, 24152, 32100-2, 32105-6, 32108-9, 32111-13, 32115, 32118-32119, 32122, 42107-26, 42129-32, and 44221; National Center for Research Resources, Grant Number: UL1RR033176; the National Center for Advancing Translational Sciences, CTSI Grant Number: UL1TR000124; the National Institute of Diabetes and Digestive and Kidney Diseases, Grant Number: DK063491

## Abstract

Red blood cell (RBC) traits provide insight into a wide range of physiological states and exhibit moderate to high heritability, making them excellent candidates for genetic studies to inform underlying biologic mechanisms. Previous RBC trait genome-wide association studies were performed primarily in European- or Asian-ancestry populations, missing opportunities to inform understanding of RBC genetic architecture in diverse populations and reduce intervals surrounding putative functional SNPs through fine-mapping. Here, we report the first fine-mapping of 6 correlated (Pearson's  $r$  range: |0.04-0.92|) RBC traits in up to 19 036 African Americans and 19 562 Hispanic/Latino participants of the Population Architecture using Genomics and Epidemiology consortium. Trans-ethnic meta-analysis of race/ethnic- and study-specific estimates for approximately 11 000 SNPs flanking 13 previously identified association signals as well as 150 000 additional array-wide SNPs was performed using inverse-variance meta-analysis after adjusting for study and clinical covariates. Approximately half of previously reported index SNP-RBC trait associations generalized to the trans-ethnic study population ( $p < 1.7 \times 10^{-4}$ ); previously unreported independent association signals within the *ABO* region reinforce the potential for multiple functional variants affecting the same locus. Trans-ethnic fine-mapping did not reveal additional signals at the *HFE* locus independent of the known functional variants. Finally, we identified a potential novel association in the Hispanic/Latino study population at the *HECTD4/RPL6* locus for RBC count ( $p = 1.9 \times 10^{-7}$ ). The identification of a previously unknown association, generalization of a large proportion of known association signals, and refinement of known association signals all exemplify the benefits of genetic studies in diverse populations.

## 1 | INTRODUCTION

Red blood cell (RBC) trait measurements are used to characterize the physiology of RBCs in both clinical and research settings, are captured by a complete blood count panel, and include the primary traits hematocrit (HCT), hemoglobin (HGB), and RBC count. Accompanying HCT, HGB, and RBC count are 3 derived traits—mean corpuscular hemoglobin (MCH), MCH concentration (MCHC), and mean corpuscular volume (MCV)—which can be used in combination with primary traits to evaluate RBC development and maintenance. Together, primary and derived RBC trait deficiencies (e.g., abnormally low HGB or excessive RBC count) cause circulatory diseases such as thalassemia, polycythemia, and genetic or nonhereditary anemias.<sup>1–5</sup> Population-specific *HBB* causal alleles for recessive diseases such as sickle-cell anemia and  $\beta$ -thalassemia have also been associated with protection against malaria and myocardial infarction, respectively, in the heterozygous state.<sup>6–8</sup> Additionally, RBC traits have been associated with stroke, cardiovascular disease (CVD) in populations with chronic kidney disease, and all-cause mortality.<sup>9–12</sup> RBC traits are therefore of substantial public health and clinical importance, yet their underlying pathophysiological mechanisms remain incompletely characterized.

As RBC traits exhibit moderate to high heritability (40%–90%), population-based genetic analysis of these phenotypes can help identify causal alleles for and inform the underlying biology of RBC-related disorders.<sup>3,13,14</sup> To date, over 80 independent association signals with one or more RBC traits have been reported, primarily in studies of European- or Asian-ancestry populations.<sup>15–24</sup> One genome-wide association study (GWAS) performed in over 16 000 African Americans identified 12 genome-wide-significant loci previously reported in European-ancestry or Japanese populations, indicating a shared role for common variants at RBC trait association signals.<sup>16</sup> However, fine-mapping of RBC trait associations identified in GWAS has had limited success narrowing broad GWAS signals to prioritize functional candidates because of large linkage disequilibrium (LD) blocks or characterizing variants that are rare or monomorphic in Europeans or Asians, as has been demonstrated for platelet count.<sup>25,26</sup> Narrowing and fine-mapping of previously identified association signals may be improved by performing analyses in ancestrally diverse populations with multi-continental admixture, including African Americans and Hispanics/Latinos.<sup>16,17,27</sup>

Here, we evaluated 32 index SNP-RBC trait associations in 11 fine-mapped MetaboChip regions, previously identified in populations of European-, Japanese-, and South Asian descent (*SPTA1*, *BCL11A*, *HFE*, *ABO*, *HK1*, *SH2B3/ATXN2*, *LIPC*, *PPCDC*, *NUTF2*, *NEUROD2*, and *TMPRSS6*) for evidence of generalization and locus refinement in African American and Hispanic/Latino participants of the Population Architecture using Genomics and Epidemiology (PAGE) consortium.<sup>28</sup> Additionally, we evaluated all SNPs genotyped on the MetaboChip for associations not previously described in any of the 6 RBC traits. These efforts will help address gaps in understanding the genetic underpinnings of RBC traits.

## 2 | MATERIALS AND METHODS

### 2.1 | Study populations

The PAGE consortium is a National Human Genome Research Institute funded effort to examine the epidemiologic architecture of genetic variants associated with human diseases and traits across diverse populations.<sup>29</sup> The following PAGE I studies contributed to this manuscript (Supplemental Materials and Methods): the Atherosclerosis Risk in Communities Study (ARIC),<sup>30</sup> the Coronary Artery Risk Disease in Young Adults study (CARDIA),<sup>31</sup> the Cardiovascular Health Study (CHS),<sup>32</sup> the Hispanic Community Health Study/Study of Latinos (HCHS/SOL),<sup>33</sup> and the Women's Health Initiative (WHI).<sup>34</sup> The Icahn Mt. Sinai School of Medicine (MSSM) contributed both African American and Hispanic/Latino study populations separately from PAGE I.<sup>29</sup> The Institutional Review Board at all participating institutions approved the study protocol and all participants gave written consent.

### 2.2 | Genotype platforms

The MetaboChip was a custom Illumina iSELECT array that contained approximately 195 000 SNPs and was designed to support large scale follow-up of putative associations for cardiovascular and metabolic traits.<sup>28</sup> Further information on genotyping and quality control is provided in the supplemental material. We defined an index SNP as a SNP reported in the GWAS catalog prior to October 1, 2016, as having a genome-wide significant association ( $5 \times 10^{-8}$ ) with at least one of the 6 RBC traits we evaluated. Index SNPs that were not directly genotyped on the MetaboChip were represented by proxies, defined as SNPs in high LD ( $r^2 \geq 0.80$ ) with the GWAS index SNP in the ancestral population in which the association was first reported. For one index SNP, rs671 (*ALDH2*), no proxy was available because this variant is specific to populations of East Asian ancestry. A total of 74% of participants were directly genotyped on the Illumina custom MetaboChip array; genotypes for the remaining participants were imputed from the Affymetrix 6.0 panel.<sup>35</sup> After QC and study-population-specific effective heterozygosity criteria were applied, 163 929 SNPs were available for analysis in African Americans and 159 467 SNPs were available for analysis in Hispanics/Latinos.

### 2.3 | Statistical analysis

We performed 4 types of analysis: (1) generalization, whereby we examined 32 index SNP associations across 6 RBC traits; (2) fine-mapping of association signals that generalized in (1); (3) testing for independent association signals for any RBC trait within one of the 11 densely genotyped regions; and (4) discovery of previously unreported associations with any RBC trait in all remaining MetaboChip regions. Only SNPs meeting an effective heterozygosity of 35 were used within each race/ethnic study population; 4814 SNPs were excluded in Hispanics/Latinos but included for African Americans, whereas 9431 were excluded for African Americans but included for Hispanics/Latinos. We examined a maximum of 8082 SNPs in African Americans and 7991 SNPs in Hispanics/Latinos (9201 SNPs total) within one of 11 regions

densely genotyped on the Illumina MetaboChip for all nondiscovery analyses.

To interpret fine-mapping results, LD was calculated in 500 kb sliding windows using PLINK (<http://pngu.mgh.harvard.edu/purcell/plink>) and African American (ARIC data), Hispanic/Latino (HCHS/SOL data), and trans-ethnic panels (randomly sampled ARIC and HCHS/SOL participant data in proportion to the racial/ethnic-specific sample population sizes).<sup>36</sup> In addition, MetaboChip LD and frequency information (but not individual-level information) was provided by the Malmö Diet and Cancer Study on 2143 control participants from a Swedish population to facilitate LD and MAF comparisons between PAGE African American and Hispanic/Latino populations and populations of European ancestry.<sup>37</sup> We used NCBI build 36 positions for regional association plots. Recombination rates were estimated from the combined HapMap phase II data.

A weighted version of generalized estimating equations (GEE; HCHS/SOL) accommodating the HCHS/SOL sampling design, relatedness, and household structure was implemented in SUGEN.<sup>38</sup> Race/ethnic-stratified linear regression was performed for all other studies (Atherosclerosis Risk in Communities [ARIC], Coronary Artery Risk Development in Young Adults [CARDIA], Cardiovascular Health Study [CHS], Icahn Mt. Sinai School of Medicine BioMe Biobank [MSSM], and WHI) using PLINK.<sup>36</sup> We evaluate the association between each quantitative RBC trait (see Supplement for RBC trait measurement methods and calculations for derived equations, Supporting Information Table S2) and a maximum of 9201 SNPs (racial/ethnic- and study-specific effective heterozygosity  $>35$ , present in more than one study in either African Americans or Hispanics/Latinos) from 11 previously identified RBC trait loci. An additive genetic model was assumed including age, sex, study center/region, and 10 ancestry principal components. Racial/ethnic-stratified and trans-ethnic study-specific association results were combined via inverse variance meta-analysis as implemented in METAL.<sup>39</sup> Genomic inflation factors were not calculated as the design of the MetaboChip purposefully emphasizes potential functional candidates, leading to expected early departure from a uniform  $p$ -value distribution.

## 2.4 | Generalization

We defined an “association signal” as a set of SNPs genotyped in a MetaboChip fine-mapped region and exhibiting linkage disequilibrium ( $r^2 \geq 0.2$  in the Malmö Diet and Cancer Study) with a previously reported genome-wide significant SNP for one or more RBC traits. For 2 or more previously reported genome-wide-significant SNPs to be considered within the same association signal in our study, those variants had to be in LD.

We next defined an “index SNP” as the most significant previously reported SNP within an association signal for each RBC trait. In instances for which multiple SNPs were published as the most significant SNP within a particular association signal for the same trait, we selected the SNP with the lowest reported  $p$ -value as the index SNP. The index SNP within an association signal may vary by trait because of differences in sample size, measurement error, and allelic

heterogeneity among other possible reasons related to genetic architecture of the traits. Therefore, we evaluated the most significant SNP reported for each association signal-trait combination rather than selecting one index SNP to evaluate in all traits for which that association signal was previously reported as genome-wide-significant, even though some of the index SNPs likely tag the same genetic association across multiple RBC traits. For example, the *SH2B3/ATXN2* association signal has been reported for multiple RBC traits with the most significant SNP differing by trait, meaning the index SNP for RBC count is rs3184504 whereas the index SNP for hematocrit is rs11065987. These 2 SNPs are in LD and likely represent the same association signal. Furthermore, while several RBC trait associations examined in this article were first reported in Japanese populations, those associations have since been generalized to European populations. European LD blocks are typically larger than for African or admixed-ancestry haplotypes, therefore we used European LD to conservatively define loci when analyzing potential independent associations in fine-mapped regions containing previously reported RBC trait GWAS associations.

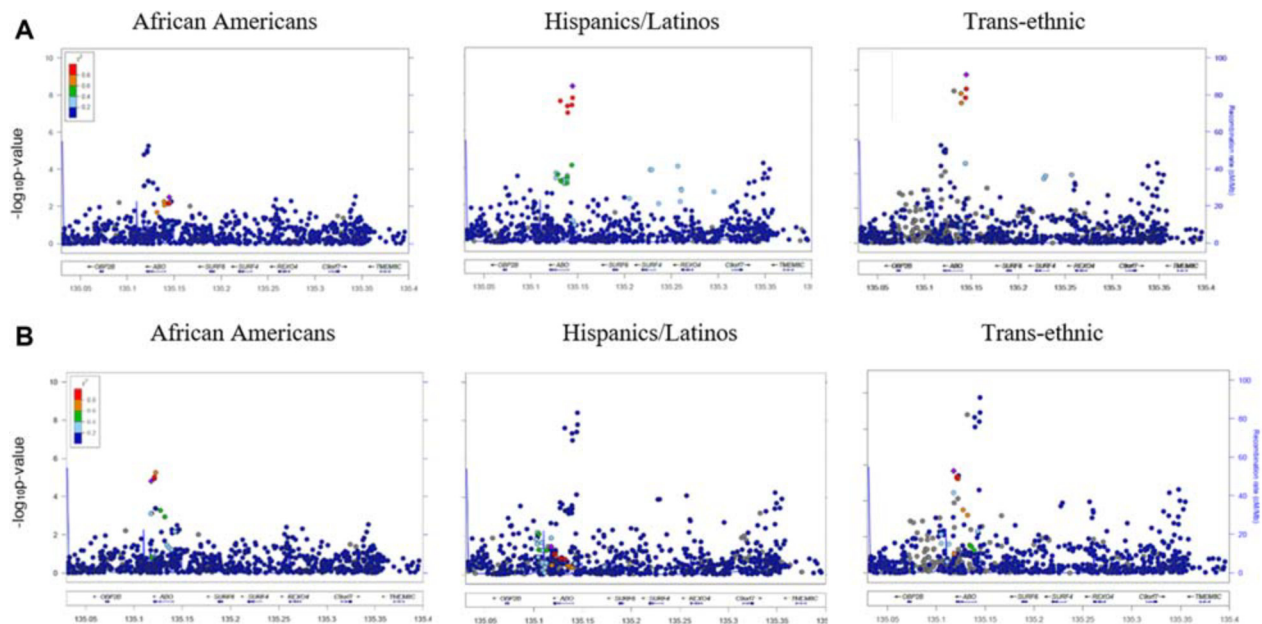
We then evaluated whether association signals identified in populations of European, Japanese, or South Asian ancestry generalized to African American and Hispanic/Latino populations. Approximately 44% of all previously reported genome-wide-significant RBC trait SNPs, but only 13% of reported association signals (defined above based on European LD, identified as of September 2016), were located in fine-mapped regions on the MetaboChip.<sup>15–24</sup> The generalization significance criterion was then defined as  $\alpha = 1.7 \times 10^{-4}$ , a Bonferroni-corrected threshold calculated using 294 tag SNPs in African Americans ( $r^2 \geq 0.80$ ; determined using African American LD from the ARIC Study) that captured all SNPs correlated with the index SNPs representing 32 index SNP-trait associations as identified in the Malmö Diet Study population.

## 2.5 | Fine-mapping generalized associations

We evaluated association-signal narrowing across ancestral backgrounds by comparing the number of SNPs in high LD with the trans-ethnic lead SNP, as well as the width of the region covered by the high-LD SNPs (Table 2, Figure 1 and Supporting Information Figure S2). LD for African Americans was calculated using ARIC study participants; LD for Hispanics/Latinos was calculated using HCHS/SOL study participants.

## 2.6 | Independent and discovery SNP identification

To identify independent SNPs influencing RBC traits, we identified all SNPs at the 11 RBC trait loci that were uncorrelated with the index SNPs ( $r^2 < 0.20$  in the Malmö Diet and Cancer Study). Sequential conditional analyses were then performed by adjusting for significant racial/ethnic-specific lead SNPs. If a statistically significant association was identified, defined as 0.05 divided by the number of SNPs in African Americans with  $MAF \geq 0.01$  that were uncorrelated with the index SNPs ( $n = 8907$ ;  $\alpha = 5.61 \times 10^{-6}$ ), the SNP was identified as independent and added to the adjustment set. Sequential conditional



**FIGURE 1** A, HGB association signal at *ABO(1)* generalizes to Hispanics/Latinos but not African Americans. B, HGB association signal at *ABO(2)* generalizes to African Americans but not Hispanics/Latinos Locus-Zoom plots representing (L-R) African American, Hispanic/Latino, and trans-ethnic meta-analysis SNP-HGB associations on respective African American, Hispanic/Latino, and trans-ethnic LD backgrounds, respectively. X-axis represents increasing chromosomal position on chromosome 9; left y-axis represents  $-\log_{10} p$ -value of each association; right y-axis represents the calculated recombination rate at each chromosomal location. Correlation of each SNP to the index SNP (purple diamond) are indicated by color as shown in the  $r$ -squared values in the box on the upper left of the left-most Locus-Zoom plots [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

analysis was repeated until no significant SNPs were identified. We evaluated all remaining SNPs for discovery in African Americans or Hispanics/Latinos using a Metachip-wide significant threshold of 0.05/155 022 (the number of SNPs available for evaluation after exclusion of SNPs evaluated for generalization), or  $\alpha = 3.23 \times 10^{-7}$ , and only considered SNPs with an effective heterozygosity  $>30$  in more than one cohort study population per race/ethnicity.

### 2.7 | Bioinformatic characterization of RBC trait loci

For each of the significant RBC trait SNPs (i.e., any lead SNP that generalized in one or both race/ethnic populations or the trans-ethnic meta-analysis; or any novel SNP identified in either race/ethnic population), all SNPs in LD ( $r^2 \geq 0.8$ ) were identified in the appropriate 1000 Genomes reference superpopulations (AFR [Africans] for African Americans and AMR [Admixed Americans] for Hispanics/Latinos) for functional annotation. Using HaploRegV2,<sup>40</sup> all variants in each LD block were characterized with putative functional roles including: conservation; promoter and/or enhancer epigenetic markers, derived from the Roadmap Epigenomics Project<sup>41</sup> and ENCODE<sup>42</sup>; DNase hypersensitive sites; and transcription factor binding motifs calculated as a library of position weight matrices.<sup>43–45</sup> Evidence of functional activity considered promoter and enhancer regions based on histone modification patterns in k562 erythroleukemia cells in the Roadmap Epigenome Project; DNase hypersensitive sites in ENCODE tissues including erythroblasts and erythroid leukemia cell lines; and transcription factor binding evidence in k562 erythroid leukemia cells. Evidence of cis-eQTL status

was performed using Blueprint for relevant blood tissues.<sup>46</sup> All functional elements that varied by cell type were restricted to RBC-relevant tissues (Supporting Information Tables S12a, S12b).

In order to evaluate the relevance of trans-ethnic PAGE lead SNPs across tissue types, we compared the eQTL status of index SNPs in both blood-relevant and other tissues (Supporting Information Table S14). We looked up significant eQTLs for each index SNP ( $p < 1E-06$ ) in whole blood in GTEx, which provides data on a wide array of tissues; and 2 blood-specific eQTL databases: the blood eQTL browser; and the NESDA NTR Conditional eQTL catalog.<sup>47–49</sup> Only GTEx tissues which showed an association with at least one SNP are reported. We further reported clinical relevance of the trans-ethnic lead SNPs as described in the literature.<sup>50,51</sup>

## 3 | RESULTS

We analyzed 6 correlated RBC traits (Pearson's correlation coefficient range:  $-0.29$  to  $0.92$  in Hispanic Community Health Study/Study of Latinos (HCHS/SOL) participants; Supporting Information Tables S1, S2) in a maximum of 19 036 African American and 19 562 Hispanic/Latino participants from 6 studies participating in the PAGE consortium (Supporting Information Table S3). Females were over-represented among both African Americans (83%) and Hispanic/Latinos (70%). The HCHS/SOL ( $n = 11\ 675$ ) and Women's Health Initiative (WHI,  $n = 17\ 363$ , of which 12 022 are African American) studies contributed the largest proportion of Hispanic/Latino (60%) and African American (63%) participants, respectively.

### 3.1 | Generalization and fine-mapping of 11 densely genotyped metabochip regions

#### 3.1.1 | Generalization

First, we examined 11 regions densely genotyped on the Illumina Metabochip and harboring one or more variants previously associated at genome-wide significance ( $p < 5 \times 10^{-8}$ ) with at least one RBC trait (Supporting Information Table S4). All but 2 of the 11 regions contained one association signal, with the *HFE* and *ABO* regions each containing 2 association signals (see Methods). Of these 13 association signals, 8 were previously associated with 2 or more RBC traits and 2 were previously associated with 4 traits, for a total of 32 index SNP-trait associations (Table 1, Supporting Information Tables S5, S6).

Seventeen of the 32 index SNP-trait associations (53%) generalized at  $p < 1.7 \times 10^{-4}$  to the trans-ethnic study population (Table 1, Supporting Information Tables S7, S8), of which 6 trans-ethnic lead SNPs were identical to the previously reported index SNP. Of the remaining 11 generalized associations, 9 trans-ethnic lead SNP  $p$ -values exceeded the index SNP  $p$ -values by at least an order of magnitude (Table 1, Supporting Information Table S10). Effect sizes for both generalized and nongeneralized association signals for index SNPs and trans-ethnic lead SNPs were consistent with previously reported estimates (Supporting Information Table S6).<sup>17</sup>

The first *HFE* association signal (index SNP: rs198846) generalized with the same trans-ethnic lead SNP (rs1799945, the functional H63D hemochromatosis variant) to all 3 previously reported traits—HGB, MCH, and MCV. Furthermore, both *ABO* association signals (Figure 1) and the *SH2B3/ATXN2* association signal generalized to all traits except RBC count. Notably, RBC count was the only trait for which none of the index SNP generalized in the trans-ethnic population; it also was the trait with the smallest sample size (46% of the maximum number of participants). Association signals for *SPTA1*, *BCL11A*, *LIPC*, *NUTF2*, *PPDC*, and *NEUROD2* did not generalize. Six nongeneralized index SNP-trait associations could not be evaluated for directional consistency because a proxy SNP was used in generalization analyses or the effect size was not reported in the initial publication (Supporting Information Table S9). For the remaining 8 nongeneralized index SNP-trait associations with sufficient information to evaluate directional consistency, 6 were directionally consistent in the trans-ethnic population. Additionally, when compared to SNP-trait associations from a previously published RBC trait GWAS, 7 of 11 PAGE lead SNPs exceeded the generalization significance threshold in 24 167 participants of the CHARGE consortium (Supporting Information Table S13).<sup>17</sup>

In race/ethnicity-specific meta-analyses (Table 2, Supporting Information Tables S6–S8), 9% ( $n = 3$ ) of index SNP-trait associations generalized to African Americans and generalization was limited to HGB. Conversely, 38% ( $n = 12$ ) of index SNP-trait associations generalized to Hispanics/Latinos, including the *SH2B3/ATXN2* association with RBC count, representing the only instance where evidence of generalization for the RBC trait was detected. Of note, HCT, HGB, and MCHC were reported for similar numbers of Hispanics/Latino and African American participants in our study population, whereas MCH, MCV, and RBC count were reported for more than twice as many Hispanics/Latinos as

African Americans, potentially contributing to the disparity in generalization by race/ethnicity.

#### 3.1.2 | Novel independent signals in 11 fine-mapped regions

We next evaluated the 11 fine-mapped regions to identify significant variants independent of published association signals by examining all SNPs that were uncorrelated with any of the index SNPs (see Methods). We identified no independent associations in previously reported regions for any of the 6 RBC traits (significance threshold:  $p < 1.3 \times 10^{-5}$ ).

#### 3.1.3 | Fine-mapping

To fine-map association signals that generalized, we then evaluated the LD structure in the trans-ethnic study population and by race/ethnicity (Figure 1 and Supporting Information Figure S2, Table 2). The median reduction in interval width was 75%, likely because of the large reduction in association signals for which the index variant was functional but fell within a large LD block in Europeans. The first *HFE* association signal showed consistent evidence of narrowing across 3 traits (113 kb decrement), with the same trans-ethnic lead SNP for all traits (the causal H63D variant rs1799945, CAF = 0.97 in African Americans, CAF = 0.88 in Hispanics/Latinos). Both *ABO* association signals (the latter of which has the determining variant for blood type B, rs8176746, as the published index SNP) fine-mapped to a limited number of SNPs in narrow LD blocks in the trans-ethnic study population, with the trans-ethnic lead SNP varying by trait. Rs855791 is a known functional coding variant in *TMPRSS6*, therefore we do not consider this signal to be narrowed.

#### 3.1.4 | Discovery

Next, all SNPs outside the 11 previously identified RBC trait-associated regions were evaluated for evidence of discovery associations ( $p < 3.03 \times 10^{-7}$ , see Methods). No SNP association exceeded Metabochip-wide significance in the total trans-ethnic study population or among African Americans for any RBC trait. However, one previously unreported association met Metabochip-wide significance in Hispanics/Latinos: rs76350043 at *HECTD4/RPL6* (chr 12q24.13,  $p = 2.5 \times 10^{-7}$ ) for RBC count (Supporting Information Table S11). This SNP was also nominally significant for HCT and HGB ( $p < 0.05$ ).

#### 3.1.5 *In silico* bioinformatics analysis

All SNPs highly correlated ( $r^2 \geq 0.9$  in relevant AFR or AMR 1000 Genomes Phase I ancestral populations) with trans-ethnic lead SNPs from generalization analysis were examined using publicly available functional prediction data for erythrocytes or erythroblastoid cell lines, as well as pathogenicity prediction (Supporting Information Tables S12a, S12b).<sup>52</sup> With the exception of the well-established *TMPRSS6* missense variant rs855791 (generalized to HCT, HGB, MCH, MCHC, and MCV), all trans-ethnic lead variants and their LD proxies were non-coding variants. Lead SNPs and their LD proxies most commonly exhibited potential regulatory effects including disruption of RBC-relevant transcription factor consensus sequences and sites exhibiting DNase I

TABLE 1 Association results examining evidence of generalization to a PAGE trans-ethnic population for 13 RCB trait association signals

Trait (N, AfAm) (N, Hisp)	Association Signal	Trans-ethnic Lead SNP <sup>b</sup>														
		Published index SNP					Trans-ethnic Population					African Americans		Hispanics/Latinos		
		Index SNP	Pop <sup>c</sup>	CA	1000G CAF Range <sup>d</sup>	SNP	CA	Beta (SE)	p-value	CAF	Index LD (r <sup>2</sup> )	p-value	CAF	Index LD (r <sup>2</sup> )	p-value	
HCT (19 036) (19 562)0	HFE (2)	rs1800562	EU	A	0.00-0.05	rs78273613	A	-0.17 (0.078)	0.03	0.98	0.67	0.21	0.86	0.39	0.08	
	ABO (1)d	rs495828	JP	T	0.13-0.22	rs635634	T	-0.21 (0.038)	5.1 × 10 <sup>-8</sup>	0.11	0.75	0.01	0.15	0.86	5.9 × 10 <sup>-7</sup>	
	HK1d	rs16926246	EU	C	0.84-1.00	rs72805692	A	-0.37 (0.061)	8.2 × 10 <sup>-10</sup>	0.98	0.08	2.7 × 10 <sup>-3</sup>	0.93	0.52	8.0 × 10 <sup>-8</sup>	
	SH2B3 / ATXN2d	rs11065987	EU	G	0.01-0.42	rs10774625	A	0.15 (0.034)	1.7 × 10 <sup>-5</sup>	0.09	0.83	3.3 × 10 <sup>-3</sup>	0.29	0.86	1.1 × 10 <sup>-3</sup>	
	TMPRSS6d	rs2143450 <sup>e</sup>	EU	A	0.43-0.70	rs855791	A	-0.22 (0.029)	8.8 × 10 <sup>-14</sup>	0.16	0.15	5.4 × 10 <sup>-3</sup>	0.44	0.73	8.1 × 10 <sup>-13</sup>	
	HFE (1)d	rs198846	EU, SA	A	0.03-0.16	rs1799945	C	-0.10 (0.017)	1.8 × 10 <sup>-9</sup>	0.97	0.24	8.6 × 10 <sup>-6</sup>	0.88	0.82	9.6 × 10 <sup>-6</sup>	
	HFE (2)	rs1800562	EU	A	0.00-0.05	rs55925606	A	-0.13 (0.036)	5.4 × 10 <sup>-4</sup>	0.99	0.90	0.05	0.98	0.84	4.3 × 10 <sup>-3</sup>	
	ABO (1)d	rs495828	JP	T	0.13-0.22	rs495828	T	-0.08 (0.012)	1.8 × 10 <sup>-10</sup>	0.14	1 <sup>e</sup>	3.2 × 10 <sup>-3</sup>	0.17	1 <sup>f</sup>	3.8 × 10 <sup>-9</sup>	
	ABO (2)d	rs7853989 <sup>e</sup>	EU, SA	T	0.06-0.20	rs10901252	C	0.06 (0.014)	2.2 × 10 <sup>-6</sup>	0.16	0.94	1.5 × 10 <sup>-5</sup>	0.07	0.97	0.04	
	HK1d	rs16926246	EU	C	0.84-1.00	rs72805692	A	-0.13 (0.021)	1.4 × 10 <sup>-10</sup>	0.98	0.08	8.1 × 10 <sup>-3</sup>	0.93	0.52	4.6 × 10 <sup>-9</sup>	
	...	rs10159477	EU, SA	A	0.00-0.14	...	...	...	...	...	0.10	...	...	...	0.56	...
	SH2B3 / ATXN2d	rs11065987	EU	A	0.58-0.99	rs10774631	A	-0.05 (0.011)	1.9 × 10 <sup>-6</sup>	0.22	0.02	1.9 × 10 <sup>-4</sup>	0.19	0.06	2.7 × 10 <sup>-3</sup>	
	...	rs3184504	EU, SA	T	0.01-0.47	...	...	...	...	...	0.03	...	...	...	0.07	...
LIPC	rs1532085	EU, SA	A	0.37-0.52	rs2414577	T	0.02 (0.009)	0.08	0.60	0.24	0.99	0.61	0.74	0.01		
TMPRSS6d	rs855791	EU	A	0.13-0.56	rs855791	A	-0.11 (0.010)	1.3 × 10 <sup>-25</sup>	0.16	1 <sup>f</sup>	2.1 × 10 <sup>-5</sup>	0.44	1 <sup>f</sup>	6.1 × 10 <sup>-23</sup>		
BCL11A	rs2540913 <sup>e</sup>	EU, SA	T	0.43-0.76	rs17027944	A	-0.0026 (0.0013)	0.05	0.30	0.34	0.08	0.20	0.17	0.01		
ABO (1)	rs495828	JP	T	0.13-0.22	rs635634	T	-0.0027 (0.0017)	0.11	0.10	0.75	0.22	0.15	0.86	0.18		
...	rs579459	EU, SA	T	0.78-0.87	...	...	...	...	...	0.74	...	...	...	0.87	...	
SH2B3 / ATXN2	rs3184504	EU, SA	T	0.01-0.47	rs10849944	T	0.0039 (0.0012)	9.9 × 10 <sup>-4</sup>	0.50	0.09	0.22	0.66	0.17	9.9 × 10 <sup>-4</sup>		
NUTF2	rs2271294	EU, SA	A	0.26-0.97	rs73612222	A	-0.0036 (0.0016)	0.02	0.87	0.02	2.8 × 10 <sup>-3</sup>	0.84	0.34	0.46		
NEUROD2	rs8182252 <sup>e</sup>	EU, SA	T	0.16-0.32	rs14050	T	-0.0026 (0.0012)	0.03	0.46	0.12	0.09	0.57	0.30	0.35		
BCL11A	rs13027161 <sup>e</sup>	EU, SA	T	0.43-0.76	rs2058703	A	0.0022 (0.0009)	0.01	0.68	0.15	0.05	0.71	0.22	0.08		
HFE (1)d	rs198846	EU, SA	A	0.03-0.16	rs1799945	C	-0.0088 (0.0014)	6.7 × 10 <sup>-10</sup>	0.97	0.24	0.09	0.89	0.82	1.8 × 10 <sup>-9</sup>		
HFE (2)d	rs1800562	NR	A	0.00-0.05	rs55925606	A	-0.0126 (0.0032)	8.8 × 10 <sup>-5</sup>	0.99	0.90	0.04	0.98	0.84	8.7 × 10 <sup>-4</sup>		

(Continues)

TABLE 1 (Continued)

Trait (N, AfAm) (N, Hisp)	Association Signal	Published index SNP			Trans-ethnic Lead SNP <sup>b</sup>				Trans-ethnic Population			African Americans			Hispanics/Latinos			
		Index SNP	Pop <sup>c</sup>	CA	1000G CAF Range <sup>d</sup>	SNP	CA	Beta (SE)	p-value	CAF	Index LD (r <sup>2</sup> )	p-value	CAF	Index LD (r <sup>2</sup> )	p-value	CAF	Index LD (r <sup>2</sup> )	p-value
...	...	rs1408272	EU	T	0.95-1.00	...	...	...	...	...	...	...	0.97	...	...	...	0.77	...
...	...	rs17342717	NR	T	0.00-0.08	...	...	...	...	...	...	...	0.73	...	...	...	0.44	...
TMPRS56d	...	rs855791	EU, SA	A	0.13-0.56	rs855791	A	-0.0096 (0.0008)	1.0 × 10 <sup>-30</sup>	0.15	1 <sup>f</sup>	2.6 × 10 <sup>-3</sup>	0.43	1 <sup>f</sup>	6.8 × 10 <sup>-27</sup>	...	...	
...	...	rs4820268	EU, NR	A	0.43-0.70	...	...	...	...	...	...	...	0.15	...	...	...	0.73	...
...	...	rs2143450	EU	A	0.43-0.70	...	...	...	...	...	...	...	0.15	...	...	...	0.73	...
MCHC (19 027) (19 553)	SPTA1	rs857684 <sup>e</sup>	EU, SA	C	0.11-0.44	rs863931	A	0.0009 (0.0003)	1.5 × 10 <sup>-3</sup>	0.33	0.04	6.7 × 10 <sup>-3</sup>	0.40	0.24	0.19	...	...	
...	...	rs857721 <sup>e</sup>	EU	C	0.11-0.44	...	...	...	...	...	...	...	0.04	...	...	...	0.24	...
ABO (2)d	...	rs8176746	JP	T	0.06-0.20	rs8176722	A	0.0018 (0.0005)	1.7 × 10 <sup>-4</sup>	0.14	0.76	2.3 × 10 <sup>-4</sup>	0.09	0.76	0.14	...	...	
TMPRS56d	...	rs855791	EU, SA	A	0.13-0.56	rs4820268	A	0.0021 (0.0003)	1.9 × 10 <sup>-11</sup>	0.73	0.15	0.04	0.54	0.73	1.4 × 10 <sup>-11</sup>	...	...	
...	...	rs4820268	NR	A	0.43-0.70	...	...	...	...	...	...	...	1 <sup>f</sup>	...	...	...	1 <sup>f</sup>	...
MCV (6397) (14 411)	SPTA1	rs3737515	EU	C	0.21-0.30	rs952094	A	-0.0008 (0.0007)	0.20	0.58	0.33	0.04	0.47	0.21	0.87	...	...	
...	...	rs243070 <sup>e</sup>	EU, SA	A	0.43-0.84	rs17402905	T	-0.0018 (0.0010)	0.06	0.86	0.30	0.37	0.84	0.20	0.12	...	...	
...	...	rs2540917 <sup>e</sup>	EU	T	0.43-0.76	...	...	...	...	...	...	...	0.30	...	...	...	0.20	...
HFE (1)d	...	rs198846	EU, SA	A	0.03-0.16	rs1799945	C	-0.0052 (0.0012)	2.6 × 10 <sup>-5</sup>	0.97	0.24	0.04	0.89	0.82	3.0 × 10 <sup>-4</sup>	...	...	
HFE (2)	...	rs1800562	NR	A	0.00-0.05	rs74662487	A	-0.0060 (0.0016)	2.7 × 10 <sup>-4</sup>	0.97	<0.01	0.02	0.94	<0.01	2.8 × 10 <sup>-3</sup>	...	...	
...	...	rs1408272	EU, SA	T	0.95-1.00	...	...	...	...	...	...	...	0.25	...	...	...	0.10	...
HK1	...	rs16926246	EU, SA	T	0.00-0.14	rs72805692	A	-0.0042 (0.0016)	9.7 × 10 <sup>-3</sup>	0.98	0.08	0.07	0.93	0.52	0.04	...	...	
PPDC	...	rs8028632	EU, SA	T	0.41-0.77	rs79269642	T	-0.0028 (0.0013)	0.02	0.94	0.02	0.37	0.90	0.11	0.04	...	...	
TMPRS56d	...	rs855791	EU, SA	A	0.13-0.56	rs855791	A	-0.0069 (0.0007)	1.0 × 10 <sup>-20</sup>	0.15	1 <sup>f</sup>	5.8 × 10 <sup>-3</sup>	0.43	1 <sup>f</sup>	2.4 × 10 <sup>-20</sup>	...	...	
...	...	rs4820268	EU, NR	A	0.43-0.70	...	...	...	...	...	...	...	0.15	...	...	...	0.73	...

<sup>a</sup>Allele frequencies obtained from HaploReg, v4.1, for 1000 Genomes Phase 3 global populations: AFR = African, AMR = American, ASN = Asian, and EUR = European. Alleles presented on the positive strand.

<sup>b</sup>Restricted to SNPs with effective heterozygosity > 30.

<sup>c</sup>Pop = published GWAS study population for the first report of the index SNP. EU = European; JP = Japanese; SA = South Asian; NR = Not Reported; Kullo, et al (2010) used electronic medical record data including patients from six potential race/ethnicity categories, but did not report the frequency of subpopulations or adjust for race/ethnicity in their linear regression analysis.

<sup>d</sup>Independent signals that generalized to the trans-ethnic PAGE population at a significance threshold of  $\alpha = 1.70 \times 10^{-4}$ .

<sup>e</sup>Index SNP not included on the MetaboChip; proxy SNP substituted (see Supporting Information Table S5).

<sup>f</sup>Index SNP and lead SNP are the same. AfAm, African American; Hisp, Hispanic; CAF, coded allele frequency; SNP, single nucleotide polymorphism.

TABLE 2 Narrowing of generalized fine-mapped RBC trait association signals using LD\* with lead SNPs

Trait	Locus	Trans-ethnic lead SNP	Trans-ethnic LD SNPs (n)	Trans-ethnic range (kb)	Malmö LD SNPs (n)	Malmö LD range (kb)	African American LD SNPs (n)	African American range (kb)	Hispanic/Latino LD SNPs (n)	Hispanic/Latino range (kb)	Locus Refinement (kb)
HGB	HFE (1)	rs1799945	4	37.2	11	150	5	46.2	5	37.2	113
MCH	HFE (1)	rs1799945	4	37.2	11	150	5	46.2	5	37.2	113
MCV	HFE (1)	rs1799945	4	37.2	11	150	5	46.2	5	37.2	113
MCH	HFE (2)	rs55925606	7	280	5	255	7	280	5	220	...
HCT <sup>a</sup>	ABO (1)	rs635634	4	5.1	5	13.1	3	13.1	5	13.1	8
HGB	ABO (1)	rs495828	2	0.7	5	13.1	1	0.7	4	13.1	12
HGB <sup>a</sup>	ABO (2)	rs10901252	3	3.4	18	79.3	5	4.6	5	14.2	76
MCHC <sup>a</sup>	ABO (2)	rs8176722	0	...	15	21.2	0	...	0	...	21
HCT	HK1	rs72805692	0	...	1	4.6	1	4.6	0	...	5
HGB	HK1	rs72805692	0	...	1	4.6	1	4.6	0	...	5
HCT <sup>a</sup>	SH2B3/ATXN2	rs10774625	10	1022	3	123	5	188	5	188	...
HGB <sup>a</sup>	SH2B3/ATXN2	rs10774631	18	173.7	41	181	18	176	19	174	7
HCT	TMPRSS6	rs855791	0	...	1	6.7	0	...	0	...	7
HGB	TMPRSS6	rs855791	0	...	1	6.7	0	...	0	...	7
MCH	TMPRSS6	rs855791	0	...	1	6.7	0	...	0	...	7
MCHC	TMPRSS6	rs4820268	0	...	1	6.7	0	...	0	...	7
MCV	TMPRSS6	rs855791	0	...	1	6.7	0	...	0	...	7

\*SNPs reported as in LD if MAF > 0.01 in the study population and  $r^2 > 0.8$  with the lead SNP.

<sup>a</sup>Previously reported association did not generalize when only study populations reporting all 6 RBC traits were considered. African American LD represented by the ARIC study population; Hispanic/Latino LD represented by the HCHS/SOL study population. PAGE LD was calculated using relevant proportions of African Americans and Hispanic/Latinos to represent trans-ethnic study population (see Methods).



activity. Several SNPs at generalized loci exhibited promise for molecular characterization, including rs198851, the HCT lead SNP in Hispanics/Latinos at the first *HFE* association signal. In k562 erythroid leukemia cells, rs198851 exhibits both DNase and enhancer activity, is an eQTL for *TRIM38*, and is located within an RNA Polymerase II ChIP-seq peak. With the exception of one association signal (*TMPRSS6* in both African Americans and Hispanics/Latinos), all independent signals contained at least one SNP with evidence for a regulatory function or cis-eQTL activity in relevant blood cell types (Supporting Information Tables S12a, S12b).

We also evaluated tissue specificity of significant eQTLs ( $p < 1E-06$ ) for each published index SNP or trans-ethnic lead SNP for all generalized association signals, as well as the putative clinical relevance of each SNP when information was available (Supporting Information Table S13).<sup>47–51</sup> eQTL results show varied evidence of tissue expression effects. Lead or index variants for *SH2B3/ATXN2* and *TMPRSS6* fine-mapped regions demonstrated no significant association with any gene expression in any tissue type. In contrast, SNPs within the second *ABO* association signal showed evidence of broad *ABO* expression across 30 tissue types. SNPs in the first *ABO* association signal show evidence of expression for multiple genes, but across fewer tissues than the second signal. Lead and index SNPs within the second *HFE* association signal were only associated with expression of genes other than *HFE* across a broad array of GTEx tissues; no lead or index SNPs for the second *HFE* association signal exhibited eQTL activity for the *HFE* gene in any tissue type. The first *HFE* association signal showed some evidence of tissue specificity in gene expression profiles—both the index SNP and trans-ethnic lead SNP exhibited cis-eQTLs for either *HFE* or several other genes—but overlap by tissue type was uncommon in this association signal.

## 4 | DISCUSSION

In this study we performed generalization, fine-mapping, and discovery analysis of 6 RBC traits in a population of over 38 000 African American and Hispanic/Latino PAGE participants. We demonstrated that genetic regions influencing RBC traits identified in European- and Asian-ancestry populations are also applicable to African American and Hispanic/Latino populations. The merits of incorporating multi-ethnic study populations in genomic studies were also displayed via locus refinement and identification of a previously unreported RBC trait association that warrants validation in future studies.

In the 11 fine-mapped regions we evaluated, over half of known index SNP-trait associations generalized to the trans-ethnic study population across all 6 RBC traits, indicating that the effects of known RBC loci are likely shared across ancestral populations. Additionally, 10 of 17 generalized associations (59%) met or exceeded the more stringent genome-wide significance threshold of  $5 \times 10^{-8}$  in the trans-ethnic study population. Although some association signals showed variation in lead SNP, the trans-ethnic lead SNPs almost always matched across traits when we restricted to participants with all traits measured (results not shown). The higher proportion of generalized associations

in Hispanics/Latinos compared to African Americans suggests that results in Hispanic/Latino populations may contribute disproportionately to the larger trans-ethnic findings. This was not surprising given the Metachip design that was enriched for European ancestral content as well as Hispanic/Latino genetic architecture, which shares more features with European-ancestry or Asian-ancestry individuals than does African American architecture.<sup>53</sup> Additionally, the SNPs designated as proxies for index SNPs discovered in European- or Japanese-ancestry individuals were almost always in much lower LD in African Americans than Hispanics/Latinos, suggesting that previously reported index SNPs are not highly effective for characterizing the genetic architecture of RBC traits in African Americans.

We also detected several instances where trans-ethnic lead SNPs showed considerably stronger evidence of association with RBC traits in our study population than previously reported GWAS index SNPs identified in primarily European or East Asian populations. By examining visualizations of generalized association signals, we further identified several cases in which the lead SNP in LD with the European index SNP was not the most significant SNP in the region, indicated differences in genetic architecture by ancestry. These findings are consistent with recent work in large trans-ethnic populations, which demonstrated considerable effect heterogeneity by genetic ancestry in GWAS index SNPs reported in studies of predominantly European populations; considerably less evidence of heterogeneity was detected when examining index SNPs identified in multi-ethnic populations.<sup>54</sup> Of particular relevance are recent demonstrations of inappropriate designation of variants which are rare in European populations as pathogenic when they are in fact common in other ancestral groups.<sup>55</sup> GWAS inclusive of diverse populations can improve the accuracy of identifying functional variants, but fine-mapping is particularly well suited to this type of exercise. In the era of precision medicine, as interest in genetic risk scores for RBC traits increases for conditions such as pregnancy, cardiovascular and neurological diseases, and mortality, studies examining generalization of reported loci to global populations will become even more important.<sup>22,56–59</sup>

Over the past decade, GWAS have identified hundreds of loci associated with RBC traits, but these findings incompletely account for the population-level variability attributable to additive genetic effects. A possible explanation for this missing heritability is that all genes expressed within RBC-relevant tissues play a role in RBC trait biology, but their identification would require near-infinite statistical power.<sup>60</sup> A recent review described the suite of genes affecting complex traits as including both “core” genes (i.e., those with tissue-specific effects crucial to one or few complex traits) and “peripheral” genes (i.e., those with broad expression profiles playing a role in many traits).<sup>61</sup> Distinguishing core from peripheral genes may inform canonical pathways for RBC traits, provide mechanistic insight into biology, and inform targets for pharmaceutical intervention.<sup>60,61</sup> Importantly, this designation occurs on a spectrum, with some genes not clearly predisposed to one class over the other. For example, hexokinase 1 (*HK1*) is highly and ubiquitously expressed, and several GWAS have identified associations within 200 kb of *HK1* for psychiatric phenotypes, autoimmune disorders, and blood metabolite levels.<sup>17,62</sup> However, *HK1* was also the only

generalized association signal in our study with evidence of a blood-specific eQTL, which localized to a narrow segment of intron 4 representing GWAS study populations of multiple ancestries. The approximately 10 kb segment contains multiple regulatory elements (e.g., DNase hypersensitivity regions and histone methylation marks), but GWAS findings to-date for this region remain restricted to RBC traits, and RBC trait index SNPs within 500 kb of this region remain restricted to this narrow genomic fragment. These results reinforce the concept that tissue-specific regulators may play an important role for individual complex traits in broadly expressed genes. In light of this information and other complex-trait GWAS findings, tissue-specific expression data and genomic information will be particularly relevant when considering candidate variants for functional studies.

Large-scale genetic evaluation of correlated traits is challenging, particularly when evaluating multiple populations and traits with variable sample sizes. Importantly, novel statistical methods scalable to GWAS that leverage correlation among phenotypes for novel locus discovery have been reported.<sup>56–59,63</sup> Such approaches seem particularly well suited for RBC traits, given evidence of a shared genetic architecture and the number of GWAS associations which have been reported in multiple RBC traits.<sup>17,22</sup> Regarding fine-mapping, extensions of correlated-phenotype methods were recently described and have similarly shown promise for reducing sets of SNPs for functional evaluation over single-trait methods. However, no studies to date have leveraged such innovations for discovery or fine-mapping of RBC traits.

Finally, we identified a potential novel association at the *HECTD4/RPL6* locus for RBC count. The *HECTD4/RPL6* locus has been previously associated with MCV (which exhibits modest correlation with RBC count and HGB), blood pressure, coronary heart disease, and multiple metabolic traits.<sup>64–68</sup> Additionally, coding mutations within several members of the ribosomal protein gene family have been causally associated with Diamond-Blackfan anemia, making an association with RBC count plausible.<sup>69,70</sup> This association signal fell within a sparsely genotyped region on the MetaboChip, and hence could not be further evaluated via fine-mapping. Evidence of association with multiple non-RBC traits should motivate larger efforts to understand whether the mechanisms underlying these associations are shared across traits or whether, for instance, tissue-specific effects relevant to each trait are represented by the same signal. Multi-ethnic fine-mapping to narrow the association signal for molecular characterization likely represents an ideal first-step, as functional variants in this region have not been described for other associated traits.

This study faced several limitations which deserve consideration. First, phenotype availability differed by study and smaller sample sizes for MCH, MCV, and RBC count likely reduced power, specifically among the African American study population. Second, 5 of 11 fine-mapped regions we evaluated (*SPTA1*, *BCL11A*, *HK1*, *LIPC*, and *TPRSS6*) also were mapped narrowly (<100 kb) with few SNPs, potentially providing insufficient coverage of African American or Hispanic/Latino genetic content to perform comprehensive fine-mapping and generalization analyses. Sparse coverage in the *SPTA1* and *BCL11A* regions could also contribute to lack of generalization at these loci, at which the respective genes have established, functional roles in RBC

development and maintenance.<sup>71–73</sup> With regard to bioinformatic characterization for functional candidate SNP evaluation, eQTL analysis—while insufficient as the sole determinant of tissue specificity—is an important component for ascertaining functional status of candidate variants. Finally, the MetaboChip design emphasized regions identified for cardiometabolic traits, so overlap with RBC-trait associations was coincidental; we therefore could not examine generalization or fine-mapping in several well established RBC trait associations, including *HBS1L/MYB*, *LUC7L/ITGF3*, *HBA1/2*, and *HBB*.<sup>17,18,21,74</sup>

## 5 | CONCLUSION

Population-based GWAS emphasize discovery, and are often the first step toward elucidating the genetic architecture underlying complex quantitative traits like RBC traits. Fine-mapping previously reported associations—particularly associations identified in genetically homogeneous populations, including European- and East Asian-ancestry populations—provides additional information about known association signals and can lead to narrowing of broad association signals to reduce the burden for bioinformatics and molecular functional analysis. Additional characterization of genetic associations contributing to population-level variability of RBC traits through large-scale sequencing and methods exploiting the correlation of RBC traits may further illuminate biological pathways for these complex quantitative traits.

## ACKNOWLEDGMENTS

(a) The Population Architecture Using Genomics and Epidemiology (PAGE) program is funded by the National Human Genome Research Institute (NHGRI), supported by U01HG004803 (CALiCo), U01HG004790 (WHI), and U01HG004801 (Coordinating Center), and their respective NHGRI ARRA supplements. The contents of this paper are solely the responsibility of the authors and do not necessarily represent the official views of the NIH. The complete list of PAGE members can be found at <http://www.pagestudy.org>.

(b) The data and materials included in this report result from a collaboration between the following studies:

Funding support for the “Epidemiology of putative genetic variants: The Women’s Health Initiative” study is provided through the NHGRI PAGE program (U01HG004790 and its NHGRI ARRA supplement). The WHI program is funded by the National Heart, Lung, and Blood Institute; NIH; and U.S. Department of Health and Human Services through contracts N01WH22110, 24152, 32100-2, 32105-6, 32108-9, 32111-13, 32115, 32118–32119, 32122, 42107-26, 42129-32, and 44221. The authors thank the WHI investigators and staff for their dedication, and the study participants for making the program possible. A full listing of WHI investigators can be found at: [http://www.whiscience.org/publications/WHI\\_investigators\\_shortlist.pdf](http://www.whiscience.org/publications/WHI_investigators_shortlist.pdf).

Funding support for the Genetic Epidemiology of Causal Variants Across the Life Course (CALiCo) program was provided through the NHGRI PAGE program (U01HG004803 and its NHGRI ARRA supplement). The following studies contributed to this manuscript and

are funded by the following agencies: The Atherosclerosis Risk in Communities (ARIC) Study is carried out as a collaborative study supported by National Heart, Lung, and Blood Institute contracts N01-HC-55015, N01-HC-55016, N01-HC-55018, N01-HC-55019, N01-HC-55020, N01-HC-55021, N01-HC-55022. The Cardiovascular Health Study (CHS) is supported by contracts HHSN268201200036C, HHSN268200800007C, N01 HC55222, N01HC85079, N01HC85080, N01HC85081, N01HC85082, N01HC85083, N01HC85086, and grants HL080295 and HL087652 from the National Heart, Lung, and Blood Institute (NHLBI), with additional contribution from the National Institute of Neurological Disorders and Stroke (NINDS). Additional support was provided by AG023629 from the National Institute on Aging (NIA). A full list of principal CHS investigators and institutions can be found at <http://www.chs-nhlbi.org/PI.htm>. CHS GWAS DNA handling and genotyping at Cedars-Sinai Medical Center was supported in part by the National Center for Research Resources, grant UL1RR033176, and is now at the National Center for Advancing Translational Sciences, CTSI grant UL1TR000124; in addition the National Institute of Diabetes and Digestive and Kidney Diseases grant DK063491 to the Southern California Diabetes Endocrinology Research Center.

Assistance with phenotype harmonization, SNP selection and annotation, data cleaning, data management, integration and dissemination, and general study coordination was provided by the PAGE Coordinating Center (U01HG004801-01 and its NHGRI ARRA supplement). The National Institutes of Mental Health also contributes to the support for the Coordinating Center.

The PAGE consortium thanks the staff and participants of all PAGE studies for their important contributions.

## ORCID

Chani Jo Hodonsky  <http://orcid.org/0000-0001-8566-5877>

## REFERENCES

- [1] Taliaferro WH, Huck JG. The inheritance of sickle-cell anaemia in man. *Genetics*. 1923;8:594–598.
- [2] Williamson GR, Crawford R. Fatal Mediterranean (Cooley's) anemia. *New Orleans Med Surg J*. 1945;98:280–284.
- [3] Whitfield JB, Martin NG. Genetic and environmental influences on the size and number of cells in the blood. *Genet Epidemiol*. 1985;2:133–144.
- [4] Lamson PD. The processes taking place in the body by which the number of erythrocytes per unit volume of blood is increased in acute experimental polycythaemia. *Proc Natl Acad Sci U S A*. 1916;2:365–369.
- [5] Neel JV, Valentine WN. Further studies on the genetics of thalassemia. *Genetics*. 1947;32:38–63.
- [6] Chami N, Lettre G. Lessons and implications from genome-wide association studies (GWAS) findings of blood cell phenotypes. *Genes (Basel)*. 2014;5(1):51–64.
- [7] Hashemi M, Shirzadi E, Talei Z, et al. Effect of heterozygous beta-thalassaemia trait on coronary atherosclerosis via coronary artery disease risk factors: a preliminary study. *Cardiovasc J Afr*. 2007;18:165–168.
- [8] Wang CH, Schilling RF. Myocardial infarction and thalassemia trait: an example of heterozygote advantage. *Am J Hematol*. 1995;49(1):73–75.
- [9] Franczuk P, Kaczorowski M, Kucharska K, et al. Could an analysis of mean corpuscular volume help to improve a risk stratification in non-anemic patients with acute myocardial infarction? *Cardiol J*. 2015;
- [10] Panwar B, Judd SE, Warnock DG, et al. Hemoglobin concentration and risk of incident stroke in community-living adults. *Stroke*. 2016;47(8):2017–2024.
- [11] Barlas RS, Honney K, Loke YK, et al. Impact of hemoglobin levels and anemia on mortality in acute stroke: analysis of UK regional registry data, systematic review, and meta-analysis. *J Am Heart Assoc*. 2016;5(8):e003019.
- [12] Solak Y, Yilmaz MI, Saglam M, et al. Mean corpuscular volume is associated with endothelial dysfunction and predicts composite cardiovascular events in patients with chronic kidney disease. *Nephrology (Carlton)*. 2013;18(11):728–735.
- [13] Evans DM, Frazer IH, Martin NG. Genetic and environmental causes of variation in basal levels of blood cells. *Twin Res*. 1999;2(04):250–257.
- [14] Wright FA, Sullivan PF, Brooks AI, et al. Heritability and genomics of gene expression in peripheral blood. *Nat Genet*. 2014;46(5):430–437.
- [15] Chambers JC, Zhang W, Li Y, et al. Genome-wide association study identifies variants in *TMPRSS6* associated with hemoglobin levels. *Nat Genet*. 2009;41(11):1170–1172.
- [16] Chen Z, Tang H, Qayyum R, et al. Genome-wide association analysis of red blood cell traits in African Americans: the COGENT Network. *Hum Mol Genet*. 2013;22(12):2529–2538.
- [17] Ganesh SK, Zakai NA, van Rooij FJ, et al. Multiple loci influence erythrocyte phenotypes in the CHARGE Consortium. *Nat Genet*. 2009;41(11):1191–1198.
- [18] Kamatani Y, Matsuda K, Okada Y, et al. Genome-wide association study of hematological and biochemical traits in a Japanese population. *Nat Genet*. 2010;42(3):210–215.
- [19] Kullo IJ, Ding K, Jouni H, Smith CY, Chute CG. Genome-wide association study of red blood cell traits using the electronic medical record. *PLoS One*. 2010;5(9):e13011A
- [20] Li J, Glessner JT, Zhang H, et al. GWAS of blood cell traits identifies novel associated loci and epistatic interactions in Caucasian and African-American children. *Hum Mol Genet*. 2013;22(7):1457–1464.
- [21] Soranzo N, Spector TD, Mangino M, et al. A genome-wide meta-analysis identifies 22 loci associated with eight hematological parameters in the HaemGen consortium. *Nat Genet*. 2009;41(11):1182–1190.
- [22] van der Harst P, Zhang W, Mateo Leach I, et al. Seventy-five genetic loci influencing the human red blood cell. *Nature*. 2012;492(7429):369–375.
- [23] Ferreira MA, Hottenga JJ, Warrington NM, et al. Sequence variants in three loci influence monocyte counts and erythrocyte volume. *Am J Hum Genet*. 2009;85(5):745–749.
- [24] Yang Q, Kathiresan S, Lin JP, Tofler GH, O'Donnell CJ. Genome-wide association and linkage analyses of hemostatic factors and hematological phenotypes in the Framingham Heart Study. *BMC Med Genet*. 2007;8 Suppl 1:S12
- [25] Gravel S, Henn BM, Gutenkunst RN, et al. Demographic history and rare allele sharing among human populations. *Proc Natl Acad Sci U S A*. 2011;108(29):11983–11988.
- [26] Schick UM, Jain D, Hodonsky CJ, et al. Genome-wide association study of platelet count identifies ancestry-specific loci in Hispanic/Latino Americans. *Am J Hum Genet*. 2016;98(2):229–242.

- [27] McCarthy MI, Hirschhorn JN. Genome-wide association studies: potential next steps on a genetic journey. *Hum Mol Genet.* 2008;17(R2):R156–R165.
- [28] Voight BF, Kang HM, Ding J, et al. The metabochip, a custom genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits. *PLoS Genet.* 2012;8(8):e1002793.
- [29] Matise TC, Ambite JL, Buyske S, et al. The Next PAGE in understanding complex traits: design for the analysis of Population Architecture Using Genetics and Epidemiology (PAGE) Study. *Am J Epidemiol.* 2011;174(7):849–859.
- [30] ARIC Investigators. The atherosclerosis risk in communities (ARIC) study: design and objectives. The ARIC investigators. *Am J Epidemiol.* 1989;129(4):687–702.
- [31] Friedman GD, Cutter GR, Donahue RP, et al. CARDIA: study design, recruitment, and some characteristics of the examined subjects. *J Clin Epidemiol.* 1988;41(11):1105–1116.
- [32] Fried LP, Borhani NO, Enright P, et al. The cardiovascular health study: design and rationale. *Ann Epidemiol.* 1991;1(3):263–276.
- [33] Daviglus ML, Talavera GA, Aviles-Santa ML, et al. Prevalence of major cardiovascular risk factors and cardiovascular diseases among Hispanic/Latino individuals of diverse backgrounds in the United States. *JAMA.* 2012;308(17):1775–1784.
- [34] WHI Study Group. Design of the Women's health initiative clinical trial and observational study. The Women's Health Initiative Study Group. *Control Clin Trials.* 1998;19(1):61–109.
- [35] Li L, Li Y, Browning SR, et al. Performance of genotype imputation for rare variants identified in exons and flanking regions of genes. *PLoS One.* 2011;6(9):e24945.
- [36] Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet.* 2007;81(3):559–575.
- [37] Minisymposium: the Malmo Diet and Cancer Study. Design, biological bank and biomarker programme. 23 October 1991, Malmo, Sweden. *J Intern Med.* 1993;233:39–79.
- [38] Lin DY, Tao R, Kalsbeek WD, et al. Genetic association analysis under complex survey sampling: the Hispanic Community Health Study/Study of Latinos. *Am J Hum Genet.* 2014;95(6):675–688.
- [39] Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol.* 2010;34(8):816–834.
- [40] Ward LD, Kellis M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.* 2012;40(D1):D930–D934.
- [41] Kundaje A, Meuleman W, Ernst J, et al. Integrative analysis of 111 reference human epigenomes. *Nature.* 2015;518(7539):317–330.
- [42] Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature.* 2012;489(7414):57–74.
- [43] Badis G, Berger MF, Philippakis AA, et al. Diversity and complexity in DNA recognition by transcription factors. *Science.* 2009;324(5935):1720–1723.
- [44] Berger MF, Badis G, Gehrke AR, et al. Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell.* 2008;133(7):1266–1276.
- [45] Berger MF, Philippakis AA, Qureshi AM, He FS, Estep PW 3rd, Bulyk ML. Compact, universal DNA microarrays to comprehensively determine transcription-factor binding site specificities. *Nat Biotechnol.* 2006;24(11):1429–1435.
- [46] Adams D, Altucci L, Antonarakis SE, et al. BLUEPRINT to decode the epigenetic signature written in blood. *Nat Biotechnol.* 2012;30(3):224–226.
- [47] Ardlie KG, Deluca DS, Segre AV, et al. Human genomics. The genotype-tissue expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science.* 2015;348(6235):648–660.
- [48] Jansen R, Hottenga JJ, Nivard MG, et al. Conditional eQTL analysis reveals allelic heterogeneity of gene expression. *Hum Mol Genet.* 2017;26(8):1444–1451.
- [49] Westra HJ, Peters MJ, Esko T, et al. Systematic identification of trans eQTLs as putative drivers of known disease associations. *Nat Genet.* 2013;45(10):1238–1243.
- [50] Landrum MJ, Lee JM, Riley GR, et al. ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucl Acids Res.* 2014;42(D1):D980–D985.
- [51] Zhou X, Maricque B, Xie M, et al. The human epigenome browser at Washington University. *Nat Methods.* 2011;8(12):989–990.
- [52] Ward LD, Kellis M. HaploReg v4: systematic mining of putative causal variants, cell types, regulators and target genes for human complex traits and disease. *Nucleic Acids Res.* 2016;44(D1):D877–D881.
- [53] Conomos MP, Reiner AP, Weir BS, Thornton TA. Model-free estimation of recent genetic relatedness. *Am J Hum Genet.* 2016;98(1):127–148.
- [54] Wojcik GL, Graff M, Nishimura KK, et al. (2017). Genetic diversity turns a new PAGE in our understanding of complex traits. (In review).
- [55] Manrai AK, Funke BH, Rehm HL, et al. Genetic misdiagnoses and the potential for health disparities. *N Engl J Med.* 2016;375(7):655–665.
- [56] Kim J, Bai Y, Pan W. An adaptive association test for multiple phenotypes with GWAS summary statistics. *Genet Epidemiol.* 2015;39(8):651–663.
- [57] Pan W, Kim J, Zhang Y, Shen X, Wei P. A powerful and adaptive association test for rare variants. *Genetics.* 2014;197(4):1081–1095.
- [58] Pasaniuc B, Price AL. Dissecting the genetics of complex traits using summary association statistics. *Nat Rev Genet.* 2017; 18(2):117–127.
- [59] Kichaev G, Roytman M, Johnson R, et al. Improved methods for multi-trait fine mapping of pleiotropic risk loci. *Bioinformatics.* 2017; 33(2):248–255.
- [60] Chakravarti A, Turner TN. Revealing rate-limiting steps in complex disease biology: the crucial importance of studying rare, extreme-phenotype families. *Bioessays.* 2016;38(6):578–586.
- [61] Boyle EA, Li YI, Pritchard JK. An expanded view of complex traits: from polygenic to omnigenic. *Cell.* 2017;169(7):1177–1186.
- [62] Rawofi L, Edwards M, Krithika S, et al. Genome-wide association study of pigmented traits (skin and iris color) in individuals of East Asian ancestry. *Peer J.* 2017;5:e3951.
- [63] Wei P, Cao Y, Zhang Y, et al. On robust association testing for quantitative traits and rare variants. *G3 (Bethesda).* 2016;6(12):3941–3950.
- [64] Astle WJ, Elding H, Jiang T, et al. The Allelic landscape of human blood cell trait variation and links to common complex disease. *Cell.* 2016;167(5):1415–1429.
- [65] Kato N, Loh M, Takeuchi F, et al. Trans-ancestry genome-wide association study identifies 12 genetic loci influencing blood pressure and implicates a role for DNA methylation. *Nat Genet.* 2015;47(11):1282–1293.
- [66] van Rooij FJ, Qayyum R, Smith AV, et al. Genome-wide trans-ethnic meta-analysis identifies seven genetic loci influencing erythrocyte traits and a role for RBPMS in erythropoiesis. *Am J Hum Genet.* 2017;100(1):51–63.
- [67] Kato N, Takeuchi F, Tabara Y, et al. Meta-analysis of genome-wide association studies identifies common variants associated with blood pressure variation in east Asians. *Nat Genet.* 2011;43(6):531–538.

- [68] Ligthart S, Vaez A, Hsu YH, et al. Bivariate genome-wide association study identifies novel pleiotropic loci for lipids and inflammation. *BMC Genomics*. 2016;17(1):443.
- [69] Cmejla R, Cmejlova J, Handrkova H, et al. Identification of mutations in the ribosomal protein L5 (RPL5) and ribosomal protein L11 (RPL11) genes in Czech patients with Diamond-Blackfan anemia. *Hum Mutat*. 2009;30(3):321–327.
- [70] Konno Y, Toki T, Tandai S, et al. Mutations in the ribosomal protein genes in Japanese patients with Diamond-Blackfan anemia. *Haematologica*. 2010;95(8):1293–1299.
- [71] Bauer DE, Orkin SH. Hemoglobin switching's surprise: the versatile transcription factor BCL11A is a master repressor of fetal hemoglobin. *Curr Opin Genet Dev*. 2015;33:62–70.
- [72] An X, Mohandas N. Disorders of red cell membrane. *Br J Haematol*. 2008;141(3):367–375.
- [73] Mankelov TJ, Satchwell TJ, Burton NM. Refined views of multi-protein complexes in the erythrocyte membrane. *Blood Cells Mol Dis*. 2012;49(1):1–10.
- [74] Hodonsky CJ, Jain D, Schick UM, et al. Genome-wide association study of red blood cell traits in Hispanics/Latinos: the Hispanic Community Health Study/Study of Latinos. *PLoS Genet*. 2017;13(4):e1006760.

#### SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

**How to cite this article:** Hodonsky CJ, Schurmann C, Schick UM, et al. Generalization and fine mapping of red blood cell trait genetic associations to multi-ethnic populations: The PAGE study. *Am J Hematol*. 2018;93:1061–1073. <https://doi.org/10.1002/ajh.25161>