

Longitudinal data reveal strong genetic and weak non-genetic components of ethnicity-dependent blood DNA methylation levels

Chris McKennan ^a, Katherine Naughton ^b, Catherine Stanhope ^b, Meyer Kattan^c, George T. O'Connor^d, Megan T. Sandel^d, Cynthia M. Visness^e, Robert A. Wood^f, Leonard B. Bacharier^g, Avraham Beigelman^g, Stephanie Lovinsky-Desir^c, Alkis Togias^h, James E. Gern ⁱ, Dan Nicolae^{b,j,*}, and Carole Ober ^{b,*}

^aDepartment of Statistics, University of Pittsburgh, Pittsburgh, PA, USA; ^bDepartment of Human Genetics, University of Chicago, Chicago, IL, USA; ^cDepartment of Pediatrics, Columbia University Medical Center, New York, NY, USA; ^dDepartment of Medicine, Boston University School of Medicine, Boston, MA, USA; ^eRho Federal Systems Division, Chapel Hill, NC, USA; ^fDepartment of Pediatrics, Johns Hopkins University Medical Center, Baltimore, MD, USA; ^gDepartment of Pediatrics, Washington University School of Medicine and St Louis Children's Hospital, St. Louis, MO, USA; ^hNational Institute of Allergy and Infectious Disease, Bethesda, MD, USA; ⁱDepartments of Pediatrics and Medicine, University of Wisconsin School of Medicine and Public Health, Madison, WI, USA; ^jDepartment of Statistics, University of Chicago, Chicago, IL, USA

ABSTRACT

Epigenetic architecture is influenced by genetic and environmental factors, but little is known about their relative contributions or longitudinal dynamics. Here, we studied DNA methylation (DNAm) at over 750,000 CpG sites in mononuclear blood cells collected at birth and age 7 from 196 children of primarily self-reported Black and Hispanic ethnicities to study race-associated DNAm patterns. We developed a novel Bayesian method for high-dimensional longitudinal data and showed that race-associated DNAm patterns at birth and age 7 are nearly identical. Additionally, we estimated that up to 51% of all self-reported race-associated CpGs had race-dependent DNAm levels that were mediated through local genotype and, quite surprisingly, found that genetic factors explained an overwhelming majority of the variation in DNAm levels at other, previously identified, environmentally-associated CpGs. These results indicate that race-associated blood DNAm patterns in particular, and blood DNAm levels in general, are primarily driven by genetic factors, and are not as sensitive to environmental exposures as previously suggested, at least during the first 7 years of life.

ARTICLE HISTORY

Received 2 April 2020
Revised 6 July 2020
Accepted 24 July 2020

KEYWORDS

DNA methylation; race/ethnicity; gene vs. environment; longitudinal epigenetics; Bayesian

Introduction

DNA methylation (DNAm) in the human genome plays a critical in regulating many cellular processes [1,2], and altered DNAm patterns have been associated with many diseases, including cancer [3], neurological disorders [4,5] and asthma [6,7], to name a few. DNAm itself reflects the contributions of genetic variation [8,9], exposure histories [10–16], and biological factors such as age [17–26], and has therefore been suggested as a mediator of the effect of these factors on disease outcomes [27,28].

Recently, results from cross-sectional studies have shown that DNAm in blood cells differ across racial and ethnic groups at birth [29,30] and later in life [31–34], suggesting that it might contribute to race/

ethnicity-associated health disparities [30,31]. Because racial and ethnic group definitions reflect both common genetic ancestries and shared diet and exposure histories [35–38], it has been postulated that race/ethnicity-associated blood DNAm patterns are an amalgam of genetic and non-genetic components, and understanding the contribution of each can help inform the relative contribution of genetic and socio-cultural diversity to variation in DNAm levels [31]. For example, a previous study [31] partitioned variation in DNAm levels into genetic and non-genetic sources, and concluded that non-genetic, socio-cultural sources had a significant impact on race/ethnicity-associated blood DNAm levels. However, that study, and all previous studies that identified race/ethnicity-associated DNAm marks, relied on cross-

CONTACT Chris McKennan  chm195@pitt.edu  Department of Statistics, University of Pittsburgh, Pittsburgh, PA, USA

*Equal contributions

 Supplemental data for this article can be accessed [here](#).

© 2020 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group. This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives License (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited, and is not altered, transformed, or built upon in any way.

sectional data and were therefore not able to assess the temporal stability of those marks. Understanding the stability of race/ethnicity-dependent DNAm present at young ages can help to determine the extent to which race/ethnicity-dependent properties of epigenetic-driven diseases can be attributed to the innate or acquired methylome [29], and identify CpGs whose DNAm is robust or sensitive to accumulated exposures. We therefore sought to fill this gap by first identifying the factors contributing to and the temporal stability of race/ethnicity-dependent blood DNAm levels, and consequently, determining the relative contributions of genetic and environmental factors to the variation in blood DNAm levels in general.

To do so, we studied global DNAm patterns at over 750,000 CpG sites on the Illumina EPIC array in cord blood mononuclear cells (CBMCs) collected at birth and in peripheral blood mononuclear cells (PBMCs) collected at 7 years of age from 196 children participating in the Urban Environment and Childhood Asthma (URECA) birth cohort study [39,40]. This cohort is part of the NIAID-funded Inner City Asthma Consortium and is comprised of children primarily of Black and Hispanic self-reported ethnicity, with a mother and/or father with a history of at least one allergic disease, and living in low socioeconomic urban areas (see O'Connor et al. [40] for details of enrolment criteria). Mothers of children in the URECA study were enrolled during pregnancy and children were followed from birth through at least 7 years of age.

The longitudinal design of the URECA study provided us with the resolution to partition genetic from non-genetic effects on race/ethnicity-associated DNAm patterns, and yielded new insight into the factors affecting DNAm patterns at CpG sites in mononuclear (immune) cells during formative

developmental years in ethnically admixed children. Using a novel statistical method that provides a general framework for analysing longitudinal genetic and epigenetic data, we show that while DNAm levels vary with chronological age, race/ethnicity-dependent DNAm patterns are overwhelmingly conserved over the first 7 years of life and that these patterns are strongly associated, and often mediated, by local genotype. Relatedly, the variation in DNAm levels at previously reported robust exposure-associated CpGs was overwhelmingly dominated by genetic rather than environmental factors in these children. Considering the results of our study and those of a recently published comprehensive review on environmental epigenetics research [41], we suggest that race/ethnicity-dependent blood DNAm levels in particular, and blood DNAm levels in general, are primarily driven by genetic factors, and are not as responsive to environmental exposures as previously suggested [31], at least during the first 7 years of life.

Results

Our study included 196 children participants in the URECA cohort who had high-quality DNA from both CBMCs and PBMCs collected at birth and age 7, respectively, available for our study [39] (see Methods). The URECA children were classified by parent- or guardian-reported race into one of the following categories: Black, $n = 147$; Hispanic, $n = 39$; White, $n = 1$; Mixed race $n = 7$, and Other, $n = 2$. A description of the study population is shown in Table 1. Genetic ancestry, assessed using principle component analysis (PCA), revealed varying proportions of African and European ancestry along PC1 (Figure 1). Because there was little separation along PC2, and no genome-wide significant correlation

Table 1. Covariates for the $n = 196$ URECA children in our study, stratified by self-reported race.

	Black	Hispanic	White	Mixed	Other
Sample Size	147	39	1	7	2
Males (%)	71 (48%)	25 (64%)	0 (0%)	4 (57%)	0 (0%)
Asthma diagnosis at age 7 (%)	38 (26%)	12 (31%)	0 (0%)	2 (29%)	0 (0%)
Gestational age at birth, in weeks (mean [range])	39.0 [34,42]	38.9 [35,41]	36.0	39.1 [37,40]	39.0 [38,40]
Sample Collection Site					
Baltimore (%)	64 (44%)	1 (3%)	1 (100%)	3 (43%)	2 (100%)
Boston (%)	17 (12%)	5 (13%)	0 (0%)	2 (29%)	0 (0%)
New York (%)	23 (16%)	32 (82%)	0 (0%)	1 (14%)	0 (0%)
St. Louis (%)	43 (29%)	1 (3%)	0 (0%)	1 (14%)	0 (0%)

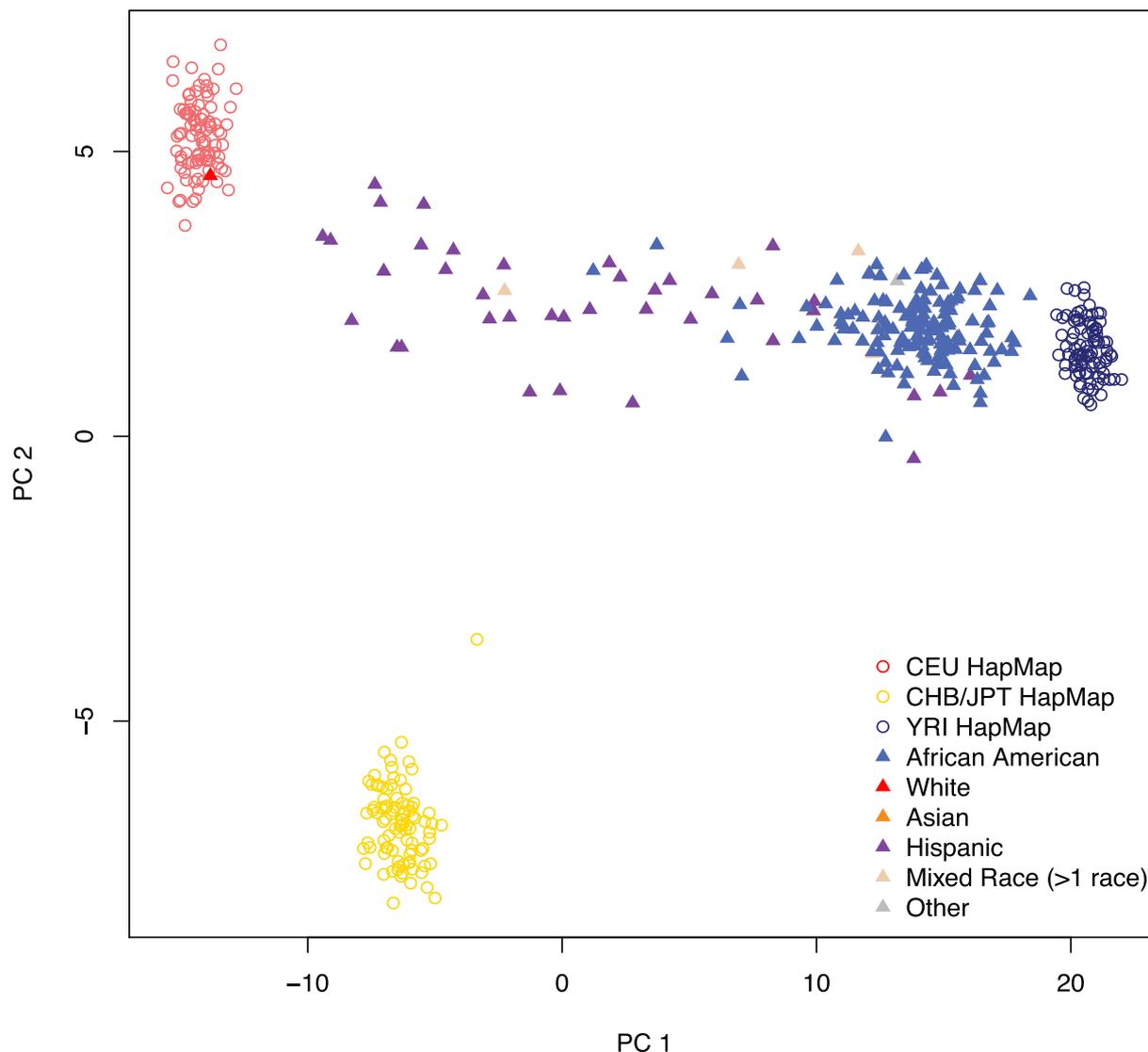


Figure 1. Estimated ancestry principal components (PCs) 1 and 2. Nearly all the variation in ancestry separates along PC1 in the URECA sample. Filled triangles represent the 196 URECA children in this study, with their self-reported race shown in different colours. Open circles are reference control samples from HapMap; red = Utah residents from northern and western Europe (CEU); yellow = east Asian (Chinese and Japanese); dark blue = Africans from Nigeria (Yoruban).

between PC2 through PC10 and DNAm levels at either age, we defined PC1 as inferred genetic ancestry. The reported races of the children are also shown in Figure 1. We included only the 186 self-reported Black and Hispanic children in subsequent analyses of reported race.

Reported race effects on DNA methylation patterns are conserved in magnitude and direction between birth and age 7

We first attempted to determine the temporal stability of reported race-associated DNAm patterns by

addressing three questions. What is the effect of reported race on DNAm levels at individual CpG sites at birth and age 7? Are the directions and magnitudes of these effects conserved from birth to age 7? Do the effects at birth and age 7 differ significantly? While these questions are important in their own right, their answers can also help determine the nature of these reported race-associated patterns. For example, race-associated DNAm levels that differ at birth and age 7 might reflect race-dependent exposure histories, while race-associated DNAm patterns that are conserved may be genetic in nature, since genetically-dependent DNAm patterns are conserved from birth to later childhood [42].

Standard hypothesis testing can be used to answer the first question but is not appropriate for answering the second or third because failure to reject the null hypothesis that the effects are equal at birth and age 7 does not imply the null hypothesis is true. Additionally, because our studies were conducted in CBMCs at birth and PBMCs at age 7, DNAm levels at birth and age 7 may differ slightly due to differences in cell composition [43]. To address these issues, we built a Bayesian model (see Model (1) in Methods) and let the data determine both the strength of the effect of reported race (based on self-report) on DNAm levels, and how similar the effects are at birth and age 7. We then answered the above three questions by defining and estimating the conserved (*con*) and discordant (*dis*) sign rates for each CpG $g = 1, \dots, 784,484$:

con_g = Posterior probability that CpG g 's reported race effects at birth and age 7 were non-zero had the same sign AND the sign was estimated correctly.

dis_g = Posterior probability that the reported race effect for CpG g was non-zero at one age and zero or in the opposite direction at the other age.

For a given posterior probability threshold, these quantities partition the reported race-associated CpGs into two groups: those whose reported race effects were non-zero and conserved from birth to age 7 and those whose reported race effects were different at birth and age 7. Detailed descriptions of our model and estimation procedure are provided in Methods and in the Supplementary Material. Supplemental Figure S1 shows how the conserved sign rate and standard P values compare.

After fitting the relevant parameters in the model to the data, we were able to estimate the fraction of CpGs with non-zero reported race effects at both ages and assign them into one of four possible bins: the two effects were completely unrelated ($\rho = 0$), moderately similar ($\rho = 1/3$), very similar ($\rho = 2/3$), or identical ($\rho = 1$). Note that if a non-trivial fraction of CpG sites had ancestry effects that were in opposite directions

at birth and age 7, they would be assigned to the first bin ($\rho = 0$). In fact, we estimated that only 0.2% of the CpGs with non-zero reported effects at both ages had unrelated or moderately similar reported race effects, whereas 30.7% fell in the very similar bin and 69.1% had identical reported race effects at birth and age 7 (Supplemental Figure S2). These data indicate that when reported race effects on DNAm levels are present (i.e., non-zero) at both birth and age 7, they tend to be very similar or exactly the same at both ages with respect to both direction and magnitude.

We then estimated the conserved and discordant sign rates for all 784,484 probes and classified a CpG as a reported race-associated CpG (RR-CpG) if its conserved or discordant sign rate was above 0.80 (i.e. $con_g \geq 0.8$ or $dis_g \geq 0.8$). At this threshold, we identified 2,162 RR-CpGs, 2,157 (99.8%) of which were conserved in sign ($con_g \geq 0.8$). Compared to self-reported Hispanic children, self-reported Black children tended to have higher DNAm levels at 1,288 (60%) of the conserved RR-CpGs ($P = 8.6 \times 10^{-38}$). This trend replicated when we substituted inferred genetic ancestry for reported race and is in accordance with previous observations [6,33], indicating individuals with more African ancestry tend to have overall more DNAm. Interestingly, there was an under enrichment of RR-CpGs in CpG islands ($P = 3.10 \times 10^{-12}$), which mirrors the observation that CpGs whose DNAm is under genetic control typically lie outside of CpG islands [44]. The fact that only 5 of the 2,162 RR-CpGs had discordant reported race effects at birth and age 7 ($dis_g \geq 0.8$) corroborates the observations made in the previous paragraph and answers the second question in the affirmative: if DNAm levels are associated with reported race at birth, the magnitude and direction of the effects are almost certainly conserved at age 7 (and vice-versa).

Inferred genetic ancestry has a larger effect on DNA methylation than does self-reported race

The observed association between self-reported race and DNAm levels may reflect differences in environmental exposures [31,33], due to associations

between race or ethnicity with socio-cultural, nutritional, and geographic exposures, among others [35–38]. In fact, a previous cross-sectional study suggested that self-reported ethnicity explained a substantial proportion of the variance of blood DNAm levels measured in Latino children of diverse ethnicities [31]. They concluded that ethnicity captured genetic, as well as the socio-cultural and environmental differences, that influence DNAm levels. If this were the case in the URECA children, the effect of inferred genetic ancestry on DNAm levels should be comparable to that of reported race. To assess this possibility in the URECA children, we repeated the analyses described above but substituted inferred genetic ancestry for reported race. This analysis revealed 8,597 inferred genetic ancestry-associated CpGs (IGA-CpGs), of which 8,579 (99.8%) were conserved in sign ($con_g \geq 0.8$). This was

significantly more than the 2,162 RR-CpGs identified in the reported race analysis above (Figure 2 (a-b)), and we show in the Supplement that this difference is robust to any differences between the powers of the reported race and inferred genetic ancestry analyses.

To further explore this finding, we examined the overlap between RR-CpGs and IGA-CpGs (Figure 2 (c)). Because reported race is an estimate of inferred genetic ancestry, there is a substantial overlap between IGA-CpGs and RR-CpGs. Contrary to the results from the previous study [31], which estimated that only 35% of their ethnicity-associated CpGs were also genetic ancestry-associated CpGs (Figure 5(a) in [31]), 66% of the RR-CpGs in our study were also IGA-CpGs, and therefore represent only a subset of the IGA-CpGs. This indicates that while IGA-CpGs include most RR-CpGs, reported race

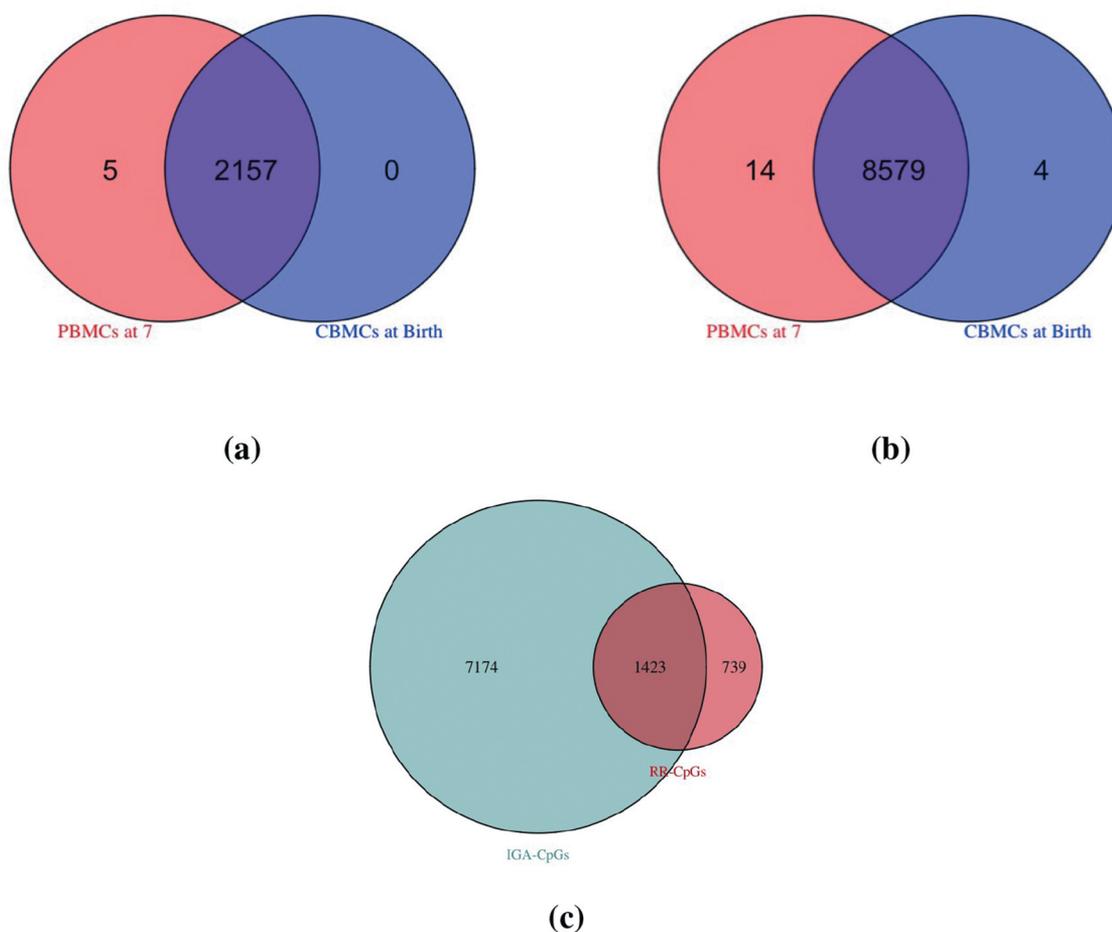


Figure 2. Overlapping ancestry CpGs at birth and at age 7. (a): self-reported race-associated CpGs (RR-CpGs) with $con_g \geq 0.8$ (violet) or $dis_g \geq 0.8$ (red or blue). A discordant RR-CpG was classified as significant at birth but not at age 7 (blue) if the marginal posterior probability that the effect was non-zero at birth was greater than that at age 7. Discordant RR-CpGs that were significant at age 7 but not at birth (red) were defined analogously. (b): The same as (a), but for inferred genetic ancestry-associated CpGs (IGA-CpGs). (c): The overlap between RR-CpGs ($con_g \geq 0.8$ or $dis_g \geq 0.8$) and IGA-CpGs ($con_g \geq 0.8$ or $dis_g \geq 0.8$).

does not capture most of the variation in DNAm levels attributable to genetic ancestry in these children.

The differences between our results and those reported in the aforementioned study may be due to the fact that sample collection site explained 80% of the variance in Mexican versus Puerto Rican ethnicity in [31], but was not accounted for in their analyses. The fact that sample collection site was associated with the DNAm levels of 865 CpGs at birth or age 7 at a 5% FDR in our study suggests that sample collection site could have confounded the relationship between ethnicity and DNAm in the previous study (see pp 7–8 in the Supplement).

The association between DNA methylation and reported race is largely genetically driven

To further address the question of whether reported race effects on DNAm levels at either birth or age 7 were primarily due to genetic variation or to environmental exposures, we used local genetic variation (within 5kb of a CpG site) and DNAm data at birth and age 7 in the 147 self-reported Black children in our study to map methylation quantitative trait loci (meQTLs). Of the 519,696 CpGs within 5kb of a SNP, 65,068 and

70,898 had at least one meQTL in CBMCs at birth and in PBMCs at age 7, respectively, at an FDR of 5%. In addition, 51% of all RR-CpGs with at least one SNP in the ± 5 kb window had at least one meQTL at birth or age 7 at an FDR of 5%, which was a significant enrichment when compared to the 17% observed for non-RR-CpGs (Figure 3(a-b)).

To provide additional evidence that local genotype mediates the effect of reported race on DNAm levels, we used logistic regression to regress the genotype of each SNP within ± 5 kb of a RR-CpG. The goal was to determine the fraction of RR-CpGs at which the observed variation was mediated through local genotype, i.e. RR-CpGs with both edges *a* and *c* in Figure 3(a). Since genotype is highly correlated with race, most SNPs will possess edge *c*. Therefore, a reasonable upper bound for this quantity is 51%, the fraction of RR-CpGs with at least one meQTL in their ± 5 kb window. To determine a lower bound, we used the results of the above-mentioned logistic regression to conservatively estimate that at least 26% of all RR-CpGs with at least one SNP in their ± 5 kb windows had both edges *a* and *c* (pp 9–11 in the Supplement). Interestingly, substituting inferred genetic ancestry for self-reported race in the above analysis yielded

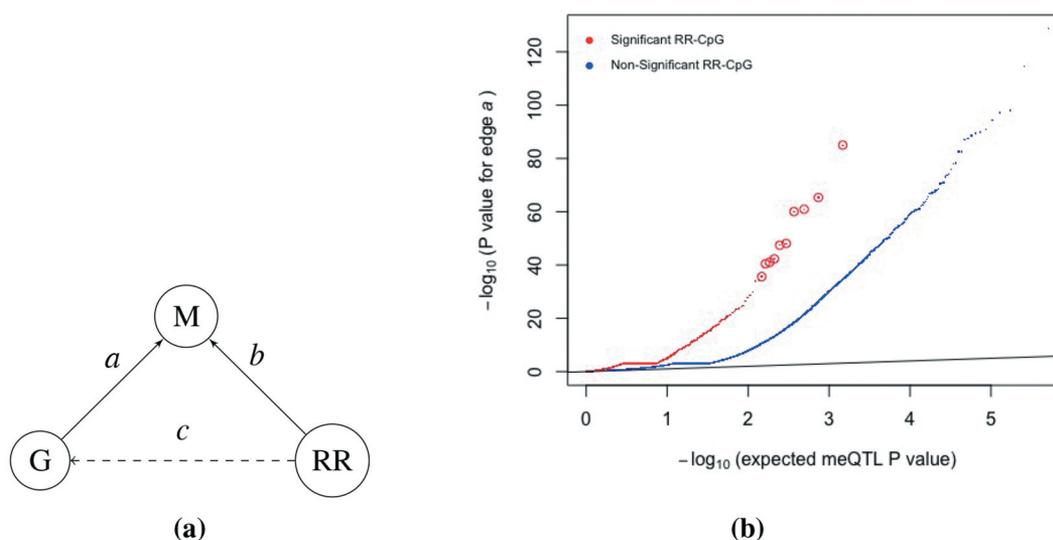


Figure 3. RR-CpGs are enriched for CpGs with meQTLs. (a) Illustration of the causal relationship between the DNAm (*m*) at a CpG site, the genotype (*g*) at the SNP within ± 5 kb of the CpG that had the smallest meQTL *P* value and self-reported race (RR). Each graph corresponds to a unique CpG. (b) Plots of the meQTL *P* value for edge *a* in CBMCs at birth, where CpGs were stratified by whether or not it was an RR-CpG ($con_g \geq 0.8$ or $dis_g \geq 0.8$). The ten enlarged red circles are just for visual aid.

nearly identical upper and lower bounds, providing evidence for local genotype mediating the effects of reported race on DNAm levels at RR-CpGs.

Genetic and biological factors explain most of the variation in blood DNA methylation levels

Given the suggested genetic nature of race/ethnicity-dependent blood cell DNAm levels, we next sought to determine the relative contributions of genetic variation, age and environmental factors on CBMC and PBMC DNAm levels in general at birth and age 7 in the URECA cohort. First, we identified 2,836 gestational age-related CpGs at birth and 16,172 age-related CpGs (CpGs whose DNAm levels changed from birth to age 7) at 5% FDRs. These two sets of CpGs were strongly enriched for CpGs used to predict gestational age in Knight et al. [21] and to predict chronological age in Horvath [18], as well as for CpGs whose blood DNAm levels changed from birth to age 5 in Pérez et al. [45] (Figure S3 in the Supplement). Moreover, the estimates of the age effects among age-related CpGs in our study showed the same direction of change as their corresponding estimated gestational age effects at birth in 97% of the 16,172 age-related CpGs. This included 14,186 gestational age-associated effects that were not significant at a 5% FDR threshold but showed the same direction of change. This concordance in direction of effect is unlikely to occur by chance (P value $< 10^{-119}$; pp 11–13 in the Supplement). Taken together with the enrichments for age-associated CpGs described above, we suggest that the majority of the changes in DNAm levels from birth to age 7 are due to ageing-related mechanisms rather than age-dependent environmental exposures.

We next attempted to determine the relative contributions of genetic and environmental factors on DNAm levels in blood. With the exception of maternal cotinine levels during pregnancy, which previously showed robust and reproducible associations with blood DNAm levels at birth [11–15] and in early childhood [10,13,16], none of the direct or indirect measures of exposures that were available in this cohort were associated with DNAm levels at either age after adjusting for

multiple testing (p 2 in the Supplement). Therefore, in order to maximize our chances of identifying environmental variation in these data, we restricted our analyses to the 6,073 maternal smoking-related CpGs identified in Joubert et al. [15], who performed a meta analysis of maternal smoking during pregnancy on 6,685 infants from 13 cohorts. In our data, DNAm levels at birth and age 7 at 505 (9.2%) and 407 (7.4%) of the 5,500 maternal smoking-related CpGs that passed QC in our study, respectively, were nominally correlated (P value ≤ 0.05) with maternal cotinine levels (enrichment P values = 7.08×10^{-34} and 6.49×10^{-8}). While this enrichment was not unexpected, we were surprised to observe that the maternal smoking-related CpGs were enriched for meQTLs (Figure 4). Additionally, there was a strong enrichment of the 8,579 conserved inferred genetic ancestry-associated CpGs among the 5,500 maternal smoking-related CpGs that passed QC in our study (fold enrichment = 2.53; P value = 6.42×10^{-33}), indicating the maternal smoking-related CpGs were enriched for genetically regulated CpGs. Furthermore, genotype at the closest SNP for over 95% of the maternal smoking-related CpGs explained a greater proportion of the variance in DNAm levels at birth than did maternal cotinine levels (Figure 4; pp 13–15 in the Supplement). These results were nearly identical for DNAm measured at age 7, and showed that genetic, and not environmental, factors are responsible for the majority of the variation in DNAm levels at even the most robust and replicated environmentally-associated CpGs in these children.

Discussion

The relationships between DNAm, chronological age, and race/ethnicity have the potential to shed light on disease aetiology and may help determine the relative genetic and environmental contributions to the observed inter-individual variability of the epigenome [17–23,29–34]. While it has previously been shown that race/ethnicity is related to DNAm in cross-sectional studies [29–34] and that statistically significant meQTLs are conserved as individuals age [42], it has yet to be shown that

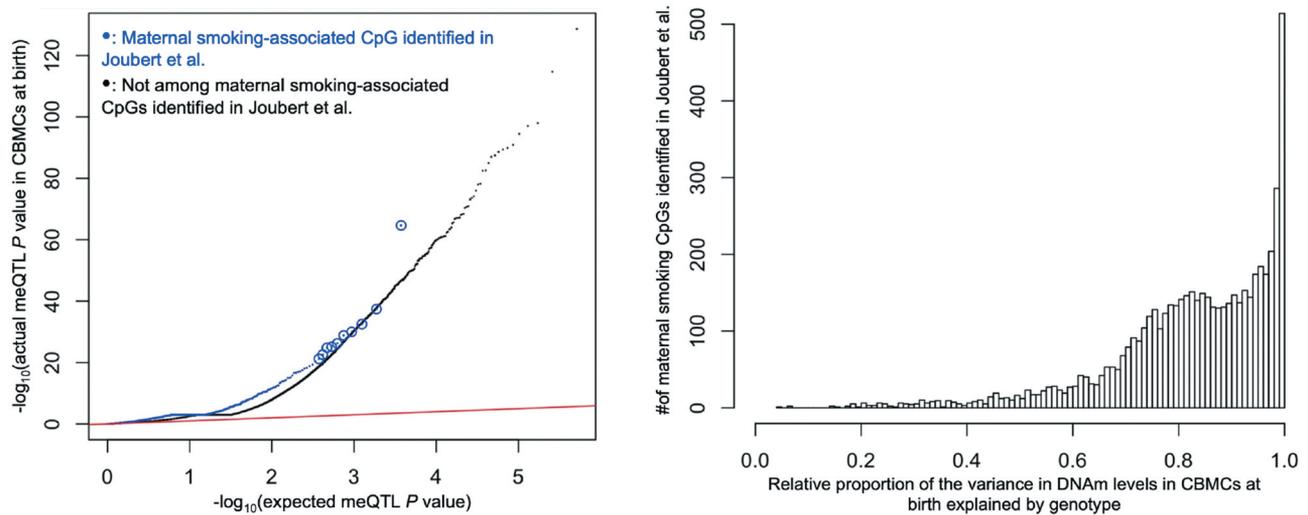


Figure 4. meQTL P value enrichment, where circled blue points are for visual aid (left), and the relative proportion of variance in DNAm levels explained by genotype (right). The x-axis of the latter was defined as the ratio of the proportion of variance in DNAm levels explained by the genotype of each CpG's closest SNP to the sum of the aforementioned genetic proportion and the proportion explained by maternal cotinine levels during pregnancy. A ratio >0.5 indicates that local genotype explained more variance than maternal cotinine levels during pregnancy.

race/ethnicity-dependent DNAm marks are conserved as children age, and relatedly, that exposure histories explain a comparatively small fraction of the variation in blood DNAm levels.

Exposure histories and other related non-genetic factors change substantially from birth to early childhood, which include changes in diet, immune profile [46], the microbiome [47] and the metabolome [48], to name a few. The putative effect of these exposures on blood DNAm [49] and the notable differences in the levels of these exposures between children of different ethnic groups [36–38] have prompted researchers to suggest that genetics only partially explain the association between ethnicity and blood DNAm levels, and that non-genetic environmental factors make a significant contribution to ethnicity-dependent blood DNAm patterns in children [29,31]. We were therefore surprised to find that self-reported race effects on DNAm were overwhelmingly conserved in both direction and magnitude from birth to age 7. This result, as well as our novel Bayesian inference paradigm used to obtain it, is important in and of itself because it provides an example of, and a general method for identifying, DNAm patterns that are conserved over time, and differentiating between environmentally responsive and temporally stable DNAm marks, which has been

highlighted as both a gap in current knowledge and a critical area of future epigenetic research [49].

While the observation that reported race effects are conserved from birth to age 7 gives credence to the hypothesis that the effects are genetic in nature, it does not rule out the possibility of environmental components or gene-environment interactions that could result in race/ethnicity-associated DNAm patterns prior to birth that persist as the child ages. It was therefore interesting to find that there was a significant under enrichment of RR-CpGs in CpG islands, which agrees with the under enrichment previously observed for CpGs under genetic control [44]. To further explore this, we showed that the RR-CpGs were enriched among CpGs with meQTLs identified in our study, indicating that DNAm levels at many of the RR-CpGs are mediated by local genotype and that much of the reported race-DNAm association could be attributed to genetic variation. Moreover, the RR-CpGs were only a small subset of IGA-CpGs in our study. Contrary to previous cross-sectional studies in infants and children [29,31], our results provide evidence for genetics accounting for an overwhelming majority of the associations between blood DNAm levels and reported race, which suggests the non-genetic contribution

to variability in blood DNAm levels may be smaller than previously thought.

There were several other notable features in these data connoting that genetic, and not environmental, factors were most responsible of the variation in blood DNAm levels in these children. The first was that although average DNAm levels of 16,172 CpGs changed significantly from birth to age 7, the direction of the change in 97% of those CpGs matched the direction of the corresponding correlation between DNAm levels and gestational age at birth. This manifest concordance in the 'epigenetic clocks' present at birth and later in life, along with the observation that the 16,172 age-related CpGs were enriched for CpGs used to predict gestational and chronological age, suggests these age-related changes are coordinated by age-related mechanisms, and not due to age-dependent environmental exposures. Second, with the exception of maternal cotinine levels during pregnancy, none of the direct or indirect measures of exposure history were associated with DNAm levels at birth or age 7. This included measures of prenatal depression and anxiety that have ostensibly been shown to be associated with cord blood DNAm patterns in other studies [50–52]. These observations are congruent with the results of a recent comprehensive review on environmental epigenetics research, which suggested that the effects of many environmental exposures on DNAm in blood are probably too small to estimate with even large sample sizes [41]. It also coincides with the rather unfortunate finding that many of the previously reported associations between exposure histories and blood DNAm are based on erroneous statistics and therefore might be spurious [53] (see pp 3–7 in the Supplement).

The third, and possibly most surprising, observation in support of strong genetically- and weak environmentally-determined blood DNAm levels were that genetic, and not maternal cotinine levels, were most responsible for the variation in DNAm levels at over 95% of the maternal smoking-associated CpGs identified in Joubert et al. [15]. This is consistent with, and significantly

extends, the results in Gonseth et al. [54], which identified genome-wide significant meQTLs for three of the top 10 most significant maternal smoking CpGs identified in the URECA5 study. It is also in line with Hannon et al. [55], which showed that genetic factors explained far more variation in the blood DNAm levels of BMI-associated CpGs than environmental factors did. One possible explanation for our observation, as demonstrated in the Gonseth et al. study, is that genotype confounds the relationship between maternal smoking and DNAm. While we did not have sufficient data to confirm this here, it remains an important area of future investigation.

Although the longitudinal features of this cohort add many strengths to our study, we must acknowledge some limitations. First, the majority of our data were derived from only two populations, self-reported Black and Hispanic children. While studying these groups makes important progress towards understanding the epigenetic architecture of underrepresented populations, it will be important to see if our conclusions replicate in other populations. Second, we only sampled DNAm through early childhood. It will be useful to assess the extent to which race/ethnicity-associated DNAm patterns persist through puberty and into adulthood.

In summary, the results of our study suggest that DNAm levels in blood cells are fairly robust to environmental exposures, including those that are associated with self-reported race. A better understanding of tissue-specific DNAm responses to environmental exposures could inform the design of future studies and provide insights into the mechanisms through which exposures and gene-environment interactions influence health and disease.

Materials and methods

Sample composition

URECA is a birth cohort study initiated in 2005 in Baltimore, Boston, New York City and St. Louis

under the NIAID-funded Inner City Asthma Consortium [39]. Pregnant women were recruited. Either they or the father of their unborn child had a history of asthma, allergic rhinitis, or eczema, and deliveries prior to 34 weeks gestation were excluded (see Gern et al. [39] for full entry criteria). Informed consent was obtained from the women at enrolment and from the parent or legal guardian of the infant after birth.

Maternal questionnaires were administered prenatally and child health questionnaires administered to a parent or caregiver every 3 months through age 7 years. Gestational age at birth and obstetric history were obtained from medical records. Additional details on study design are described in Gern et al. [39]. Frozen paired cord blood mononuclear cells (CBMCs) and peripheral blood mononuclear cells (PBMCs) at age 7, were available for 196 of the 560 URECA children after completing other studies. After QC, DNAm data were available for 194 children at birth, 195 children at age 7, and 193 children at both time points; genotype data were available in 193 children. The sample size for each analysis is given in Table 2.

Maternal cotinine levels were measured in the cord blood plasma at birth, and we categorized mothers as smokers ($\geq 10\text{ng/mL}$; $n = 31$) or non-smokers ($< 10\text{ng/mL}$; $n = 150$), where cotinine levels were missing in 15 mothers. The 10 ng/mL threshold was the same as that used in Joubert et al. [15] to define a pregnant mother with a sustained smoking habit, where 147/150 (98%) of the non-smokers in our data

had cotinine levels below 2 ng/mL, the detection limit of the assay.

DNA methylation

DNA for methylation studies was extracted from thawed CBMCs and PBMCs using the Qiagen AllPrep kit (QIAGEN, Valencia, CA). Genome-wide DNA methylation was assessed using the Illumina Infinium MethylationEPIC BeadChip (Illumina, San Diego, CA) at the University of Chicago Functional Genomics Facility (UC-FGF). Birth and 7-year samples from the same child were assayed on the same chip and the data were processed using Minfi [56]; Infinium type I and type II probe bias were corrected using SWAN [57]. Raw probe values were corrected for colour imbalance and background by control normalization. Three out of the 392 samples (two at birth and one at age 7) were removed as outliers following normalization. We removed 82,352 probes that mapped either to the sex chromosomes or to more than one location in a bisulphite-converted genome, had detection P values greater than 0.01% in 25% or more of the samples, or overlapped with known SNPs with minor allele frequency of at least 5% in African, American, or European populations. After processing, 784,484 probes were retained and M -values were used for all downstream analyses, which were computed as $\log_2(\text{methylated intensity} + 100) - \log_2(\text{unmethylated intensity} + 100)$. The offset of 100 was recommended in Du et al. [58].

Genotyping

DNA from the 196 URECA children was genotyped with the Illumina Infinium CoreExome+Custom array. Of the 532,992 autosomal SNPs on the array, 531,755 passed Quality control (QC) (excluding SNPs with call rate $< 95\%$, Hardy-Weinberg P values $< 10^{-5}$, and heterozygosity outliers). We conducted all analyses in 293,696 autosomal SNPs with a minor allele frequency $\geq 5\%$. Genotypes for three children failed QC and were excluded from subsequent analysis that involved genotypes, including methylation quantitative locus (meQTL) mapping, inferred genetic ancestry, or used genetic ancestry PC1 as a covariate. These three children were included in all other analyses.

Table 2. Sample size and composition for each analysis.

	Black	Hispanic	White	Mixed	Other
Inferred genetic ancestry, paired samples	143	37	0	0	0
Self-reported race, paired samples	145	38	0	0	0
Age (birth to age 7), paired samples	143	37	1	7	2
Gestational age at birth	144	37	1	7	2
meQTLs at birth	144	0	0	0	0
meQTLs at age 7	144	0	0	0	0
Maternal cotinine levels at birth*	132	38	1	6	2
Maternal cotinine levels at age 7*	134	37	1	6	2

*15 of the mothers did not have cord blood plasma cotinine measurements.

Estimating inferred genetic ancestry

Ancestral principal component analysis (PCA) was performed using a set of 801 ancestry informative markers (AIMs) from Tandon et al. [59] that were genotyped in both the URECA children and in HapMap [60] release 23.

Univariate statistical methods

To determine the effect of gestational age and maternal cotinine levels (smoker vs. non-smokers) on DNAm levels in CBMCs at birth or PBMCs at age 7, we used standard linear regression models with the child's gender, sample collection site, inferred genetic ancestry and methylation plate number as covariates in our model. We controlled for gestational age in the maternal cotinine analysis. We also estimated cell composition and other unobserved confounding factors using a method described in McKennan et al. [61]. We then computed P values for each CpG site and used q -values [62] to control the false discovery rate at a nominal level. We took the same approach to determine CpGs whose DNAm changed from birth to age 7, except the response was measured as the difference in DNAm at birth and age 7. In this analysis, we included the child's gender, gestational age at birth, inferred genetic ancestry and sample collection site as covariates. Because all paired samples were on the same plate, we did not include plate number as a covariate in this analysis. We also estimated unobserved factors that influence differences in DNAm at birth and age 7 using McKennan et al. [61] and included these latent factors in our linear model.

Joint modelling of DNA methylation at birth and age 7

We used data from the self-reported Hispanic and Black individuals with DNAm measured at both time points to analyse the effect of ancestry on DNAm levels at CpGs $g = 1, \dots, p = 784, 484$ using the following model:

$$\begin{aligned} y_g &= \begin{pmatrix} y_g^{(0)} \\ y_g^{(7)} \end{pmatrix} \\ &= \begin{pmatrix} X\beta_g^{(0)} \\ X\beta_g^{(7)} \end{pmatrix} + Zy_g + C\ell_g + e_g, \end{aligned} \quad (1a)$$

$$\begin{aligned} \begin{pmatrix} \beta_g^{(0)} \\ \beta_g^{(7)} \end{pmatrix} &= (\sigma_g^2 + \delta_g^2)^{-1/2} \begin{pmatrix} \beta_g^{(0)} \\ \beta_g^{(7)} \end{pmatrix} \sim \pi_{(0,0)} \delta_{(0,0)} \\ &+ \sum_{k=1}^K \pi_{(1,0)}^{(k)} \begin{pmatrix} N_1(0, \tau_k^2) \\ \delta_0 \end{pmatrix} \\ &+ \sum_{k=1}^K \pi_{(0,1)}^{(k)} \begin{pmatrix} \delta_0 \\ N_1(0, \tau_k^2) \end{pmatrix} \\ &+ \sum_{s=1}^S \sum_{k=1}^K \pi_{(1,1)}^{(k,s)} N_2 \left(0, \tau_k^2 \begin{pmatrix} 1 & \rho_s \\ \rho_s & 1 \end{pmatrix} \right), \end{aligned} \quad (1b)$$

$$\begin{aligned} e_g &\sim N_{2n} \left(0, \sigma_g^2 I_{2n} + \delta_g^2 B \right), B_{ij} \\ &= 1 \left\{ \begin{array}{l} \text{samples } i \text{ and } j \text{ are} \\ \text{from the same child} \end{array} \right\}, \end{aligned} \quad (1c)$$

where δ_0 and $\delta_{(0,0)}$ are the point masses at $0 \in \mathbb{R}$ and $(0, 0) \in \mathbb{R}^2$. The vector $y_g^{(a)} \in \mathbb{R}^n$ contained the DNAm levels at CpG g at age a , $X \in \mathbb{R}^n$ contained each child's inferred genetic ancestry or self-reported race and $\beta_g^{(a)}$ was the effect due to ancestry at age a . X was standardized to have variance 1 when X was inferred genetic ancestry. The nuisance covariates Z contained an intercept for the cord blood and PBMC samples, sample collection site, gender, gestational age at birth and plate number. Since gestational age was only correlated with cord blood DNAm, we assumed the effect of gestational age on DNAm at age 7 was zero for all CpG sites. We estimated the unobserved covariates C with McKennan et al. [63], which accounts for the correlation between samples from the same child.

The entries of the weight vector $\pi = (\pi_{(0,0)}, \pi_{(1,0)}^{(1)}, \dots, \pi_{(1,0)}^{(K)}, \pi_{(0,1)}^{(1)}, \dots, \pi_{(0,1)}^{(K)}, \pi_{(1,1)}^{(1,1)}, \dots, \pi_{(1,1)}^{(S,K)})^T$ sum to 1, where we set $K = 5$ and $S = 4$. Similar to Flutre et al. [64] and Stephens [65], we specified a grid of correlation coefficients $\rho_s \in \{0, 1/3, 2/3, 1\}$ and a dense grid of effect sizes $\tau_k \in \{0.05, 0.1, 0.15, 0.20, 0.25\}$ when X was inferred genetic ancestry and $\tau_k \in \{0.1, 0.15, 0.225, 0.3, 0.375\}$ when X was reported

race. We set τ_4 by first performing a univariate analysis and then estimating the variance of the effect sizes for CpGs with q -values ≤ 0.05 , and τ_1 was such that if $b_g^{(a)} \sim N_1(0, \tau_1^2)$, the expected number of CpGs significant at the Bonferroni threshold $0.05/p$ in a univariate analysis would be smaller than 1 for $a = 0, 7$. The proportion of CpGs with non-zero reported race effects at both ages that fell in bin $s = 1, \dots, 4$ was defined as $\sum_{k=2}^K \pi_{(1,1)}^{(k,s)}$, where

we ignored the proportion when $k = 1$, because τ_1 was too small to differentiate from zero. The estimated proportion of CpGs in the $\rho_s = 2/3$ or $\rho_s = 1$ bins was still over 98% when we included τ_1 .

To fit the model, we first regressed out Z and the estimated C from both y_g and $X \oplus X$ and used the residuals in the downstream analysis. We estimated σ_g^2 and δ_g^2 for each $g = 1, \dots, p$ with restricted maximum likelihood (REML) and followed Stephens [65] and estimated π by empirical Bayes via expectation maximization. Supplemental Figures S2 and S4 plot the estimate for π in the reported race analysis. We then defined con_g and dis_g for each CpG $g = 1, \dots, p$ as

$$con_g = \hat{\mathbb{P}} \left\{ \beta_g^{(0)}, \beta_g^{(7)} > 0 \mid \mathbf{y}_g, \pi, \sigma_g^2, \delta_g^2 \right\} \\ \vee \hat{\mathbb{P}} \left\{ \beta_g^{(0)}, \beta_g^{(7)} < 0 \mid \mathbf{y}_g, \pi, \sigma_g^2, \delta_g^2 \right\}$$

$$dis_g = \hat{\mathbb{P}} \left[\left\{ \beta_g^{(0)} > 0, \beta_g^{(7)} \leq 0 \right\} \cup \left\{ \beta_g^{(0)} < 0, \beta_g^{(7)} \geq 0 \right\} \right. \\ \left. \cup \left\{ \beta_g^{(0)} \geq 0, \beta_g^{(7)} < 0 \right\} \cup \left\{ \beta_g^{(0)} \leq 0, \beta_g^{(7)} > 0 \right\} \right. \\ \left. \mid \mathbf{y}_g, \sigma_g^2, \delta_g^2, \pi \right].$$

Determining meQTLs

We performed meQTL mapping in the 145 genotyped, self-reported Black children using the set of 269,622 SNPs with 100% genotype call rate in this subset. We restricted ourselves to this subset of samples to minimize heterogeneity in effect sizes. To identify CpG-SNP pairs, we considered SNPs within 5kb of each CpG, as this region has been previously shown to contain the majority of genetic variability in

DNAm [8] and is small enough to mitigate the multiple testing burden, and computed a P value for the effect of the genotype at a single SNP on DNAm at the corresponding CpG with ordinary least squares. We then defined the meQTL for each CpG site as the SNP with the lowest P value. In addition to genotype, we included inferred genetic ancestry (i.e., ancestry PC1), gestational age at birth, gender, sample collection site and methylation plate number in the linear model, along with the first nine principal components of the residual DNAm data matrix after regressing out the intercept and the five additional covariates. We then tested the null hypothesis that a CpG did not have an meQTL in the 10kb region by using the minimum marginal P value in the region as the test statistic and computed its significance via bootstrap. We used q -values to control the false discovery rate.

Ethical statement

We used de-identified single nucleotide polymorphism, DNA methylation and phenotype data from samples taken from human subjects as part of the Urban Environment and Childhood Asthma study. The WIRB approved human samples to be used in the Urban Environment and Childhood Asthma study (WIRB project number: 20,142,570).

Disclosure statement

No potential conflict of interest was reported by the authors.

Funding

This work was supported by the National Center for Research Resources [RR00052, M01RR00533, 1UL1RR025771, M01RR00071, 1UL1RR024156, UL1TR000040, UL1TR001079 and 5UL1RR024992-02]; National Heart, Lung, and Blood Institute [R01 HL129735]; National Heart, Lung, and Blood Institute [R01 HL122712]; National Heart, Lung, and Blood Institute [P01 HL070831]; National Institute of Allergy and Infectious Diseases [NO1-AI-25496, NO1-AI-25482, HHSN272200900052C, HHSN272201000052I, 1UM1AI114271-01 and UM2AI117870]; National Institute of Allergy and Infectious Diseases [U19 AI106683].

ORCID

Chris McKenna  <http://orcid.org/0000-0003-2096-8257>
Katherine Naughton  <http://orcid.org/0000-0002-0378-2437>

Catherine Stanhope  <http://orcid.org/0000-0001-8087-1970>
 James E. Gern  <http://orcid.org/0000-0002-6667-4708>
 Carole Ober  <http://orcid.org/0000-0003-4626-9809>

References

- [1] Bird A. DNA methylation patterns and epigenetic memory. *Genes Dev.* 2002;16:6–21.
- [2] Smith ZD, Meir A. DNA methylation: roles in mammalian development. *Nat Rev Genet.* 2013;14(3):204–220.
- [3] Baylin SB, Jones PA. Epigenetic determinants of cancer. *Cold Spring Harb Perspect Biol.* 2016 Sep;8(9):a019505.
- [4] Ladd-Acosta C, Hansen KD, Briem E, et al. Common DNA methylation alterations in multiple brain regions in autism. *Mol Psychiatry.* 2014 Sep;19:862–871.
- [5] Azkargorta MA, Urdániz-Casado A, Sánchez-Ruiz de Gardoaj, et al. DNA methylation alterations in Alzheimer's disease. *Clin Epigenet.* 2019;11:91.
- [6] Chan MA, Ciaccio CE, Gigliotti NM, et al. DNA methylation levels associated with race and childhood asthma severity. *J Asthma.* 2017 Sep;54(8):825–832.
- [7] Nicodemus-Johnson J, Myers RA, Sakabe NJ, et al. DNA methylation in lung cells is associated with asthma endotypes and genetic risk. *JCI Insight.* 2016 Dec;1(20).
- [8] Bell JT, Pai AA, Pickrell JK, et al. DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biol.* 2011;12(1):R10.
- [9] Smith AK, Kilaru V, Kocak M, et al. Methylation quantitative trait loci (meQTLs) are consistently detected across ancestry, developmental stage, and tissue type. *BMC Genomics.* 2014;15(1):145.
- [10] Breton CV, Siegmund KD, Joubert BR, et al. Prenatal tobacco smoke exposure is associated with childhood DNA CpG methylation. *Plos One.* 2014 Jun;9(6):e99716.
- [11] Markunas CA, Xu Z, Harlid S, et al. Identification of DNA methylation changes in newborns related to maternal smoking during pregnancy. *Environ Health Perspect.* 2014 Oct;122(10):1147–1153.
- [12] Ivorra C, Fraga MF, Bayón GF, et al. DNA methylation patterns in newborns exposed to tobacco in utero. *J Transl Med.* 2015 Jan;13(1):25.
- [13] Richmond RC, Simpkin AJ, Woodward G, et al. Prenatal exposure to maternal smoking and offspring DNA methylation across the lifecourse: findings from the Avon Longitudinal Study of Parents and Children (ALSPAC). *Hum Mol Genet.* 2015 Apr;24(8):2201–2217.
- [14] Joubert BR, Håberg SE, Nilsen RM, et al. 450K epigenome-wide scan identifies differential DNA methylation in newborns related to maternal smoking during pregnancy. *Environ Health Perspect.* 2012 Oct;120(10):1425–1431.
- [15] Joubert BR, Felix J, Yousefi P, et al. DNA methylation in newborns and maternal smoking in pregnancy: genome-wide consortium meta-analysis. *Am J Hum Genet.* 2016;98(4):680–696.
- [16] Rzehak P, Saffery R, Reischl E, et al. Maternal smoking during pregnancy and DNA-methylation in children at age 5.5 years: epigenome-wide-analysis in the European Childhood Obesity Project (CHOP)-study. *Plos One.* 2016 May;11(5):e0155554.
- [17] Bocklandt S, Lin W, Sehl ME, et al. Epigenetic predictor of age. *Plos One.* 2011 Jun;6(6):e14821.
- [18] Horvath S. DNA methylation age of human tissues and cell types. *Genome Biol.* 2013;14(10):3156.
- [19] Horvath S, Erhart W, Brosch M, et al. Obesity accelerates epigenetic aging of human liver. *Proc Nat Acad Sci.* 2014;111(43):15538–15543.
- [20] Johnson AA, Akman K, Calimport SRG, et al. The role of DNA methylation in aging, rejuvenation, and age-related disease. *Rejuvenation Res.* 2012;15(5):483–494.
- [21] Knight AK, Jeffrey MC, Christiane T, et al. An epigenetic clock for gestational age at birth based on blood methylation data. *Genome Biol.* 2016;17(1):206.
- [22] Levine ME, Crimmins EM. Evidence of accelerated aging among African Americans and its implications for mortality. *Soc Sci Med.* 2014;118:27–32.
- [23] Marioni RE, Shah S, McRae AF, et al. DNA methylation age of blood predicts all-cause mortality in later life. *Genome Biol.* 2015;16(1):25.
- [24] Parets SE, Conneely KN, Kilaru V, et al. Fetal DNA methylation associates with early spontaneous preterm birth and gestational age. *Plos One.* 2013 Jun;8(6):e67489.
- [25] Schroeder JW, Conneely KN, Cubells JF, et al. Neonatal DNA methylation patterns associate with gestational age. *Epigenetics.* 2011 Dec;6(12):1498–1504.
- [26] Davey Smith G, Tilling K, Suderman M, et al. Longitudinal analysis of DNA methylation associated with birth weight and gestational age. *Hum Mol Genet.* 2015 Apr;24(13):3752–3763.
- [27] Wu D, Yang H, Winham SJ, et al. Mediation analysis of alcohol consumption, DNA methylation, and epithelial ovarian cancer. *J Hum Genet.* 2018;63(3):339–348.
- [28] Huang JV, Cardenas A, Colicino E, et al. DNA methylation in blood as a mediator of the association of midchildhood body mass index with cardio-metabolic risk score in early adolescence. *Epigenetics.* 2018;13(10–11):1072–1087.
- [29] Adkins RM, Krushkal J, Tylavsky FA, et al. Racial differences in genespecific DNA methylation levels are present at birth. *Birth Defects Res A Clin Mol Teratol.* 2011;91(8):728–736.
- [30] Mozhui K, Smith AK, Tylavsky FA. Ancestry dependent DNA methylation and influence of maternal nutrition. *Plos One.* 2015 Mar;10(3):e0118466.
- [31] Galanter JM, Gignoux CR, Oh SS, et al. Differential methylation between ethnic sub-groups reflects the effect of genetic ancestry and environmental exposures. *eLife.* 2017 Jan;6:e20532.
- [32] Heyn H, Sebastian M, Irene HH, et al. DNA methylation contributes to natural human variation. *Genome Res.* 2013 Sep;23(9):1363–1372.

- [33] Moen EL, Zhang X, Mu W, et al. Genome-wide variation of cytosine modifications between European and African populations and the implications for complex traits. *Genetics*. 2013 Aug;194(4):987.
- [34] Rahmani E, Shenhav L, Schweiger R, et al. Genome-wide methylation data mirror ancestry information. *Epigenetics Chromatin*. 2017;10(1):1.
- [35] Nguyen AB, Moser R, Chou WY. Race and health profiles in the United States: an examination of the social gradient through the 2009 CHIS adult survey. *Public Health*. 2014;128(12):1076–1086.
- [36] Salvo D, Frediani JK, Ziegler TR, et al. Food group intake patterns and nutrient intake vary across low-income hispanic and African American preschool children in Atlanta: a cross sectional study. *Nutr J*. 2012;11(1):62.
- [37] Skala K, Chuang R-J, Evans A, et al. Ethnic differences in the home food environment and parental food practices among families of low income Hispanic and African-American preschoolers. *J Immigr Minor Health*. 2012 Dec;14(6):1014–1022.
- [38] Chakraborty J, Zandbergen PA. Children at risk: measuring racial/ethnic disparities in potential exposure to air pollution at school and home. *J Epidemiol Community Health*. 2007 Dec;61(12):1074–1079.
- [39] Gern JE, Visness CM, Gergen PJ, et al. The Urban Environment and Childhood Asthma (URECA) birth cohort study: design, methods, and study population. *BMC Pulm Med*. 2009;9(1):17.
- [40] O'Connor GT, Lynch SV, Bloomberg GR, et al. Early-life home environment and risk of asthma among inner-city children. *J Allergy Clin Immunol*. 2018;141(4):1468–1475.
- [41] Breton CV, Marsit CJ, Faustman E, et al. Small-magnitude effect sizes in epigenetic end points are important in children's environmental health studies: the children's environmental health and disease prevention research center's epigenetics working group. *Environ Health Perspect*. 2017 Apr;125(4):511–526.
- [42] Gaunt TR, Shihab HA, Hemani G, et al. Systematic identification of genetic influences on methylation across the human life course. *Genome Biol*. 2016;17(1):61.
- [43] Fu J, Wolfs MGM, Deelen P, et al. Unraveling the regulatory mechanisms underlying tissue-dependent genetic variation of gene expression. *PLoS Genet*. 2012 Jan;8(1):e1002431.
- [44] Lin D, Chen J, Perrone-Bizzozero N, et al. Characterization of cross-tissue genetic-epigenetic effects and their patterns in schizophrenia. *Genome Med*. 2018 Feb;10(1):13.
- [45] Pérez RF, Santamarina P, Tejedor JR, et al. Longitudinal genome-wide DNA methylation analysis uncovers persistent early-life DNA methylation changes. *J Transl Med*. 2019 Jan;17(1):15.
- [46] Olin A, Henckel E, Chen Y, et al. Stereotypic immune system development in newborn children. *Cell*. 2018 Aug;174(5):1277–1292.
- [47] Mohammadkhah AI, Simpson EB, Patterson SG, et al. Development of the gut microbiome in children, and lifetime implications for obesity and cardiometabolic disease. *Children (Basel)*. 2018 Nov;5(12):160.
- [48] Chiu C-Y, Yeh K-W, Lin G, et al. Metabolomics reveals dynamic metabolic changes associated with age in early childhood. *Plos One*. 2016 Feb;11(2):e0149823.
- [49] Martin EM, Fry RC. Environmental influences on the epigenome: exposure-associated DNA methylation in human populations. *Annu Rev Public Health*. 2018;39(1):309–333.
- [50] Hompes T, Izzi B, Gellens E, et al. Investigating the influence of maternal cortisol and emotional state during pregnancy on the DNA methylation status of the glucocorticoid receptor gene (NR3C1) promoter region in cord blood. *J Psychiatr Res*. 2013;47(7):880–891.
- [51] Non AL, Binder AM, Kubzansky LD, et al. Genome-wide DNA methylation in neonates exposed to maternal depression, anxiety, or SSRI medication during pregnancy. *Epigenetics*. 2014 Jul;9(7):964–972.
- [52] Oberlander TF, Weinberg J, Papsdorf M, et al. Prenatal exposure to maternal depression, neonatal methylation of human glucocorticoid receptor gene (NR3C1) and infant cortisol stress responses. *Epigenetics*. 2008 Mar;3(2):97–106.
- [53] Mansell T, Vuillermin P, Ponsonby A-L, et al. Maternal mental well-being during pregnancy and glucocorticoid receptor gene promoter methylation in the neonate. *Dev Psychopathol*. 2016;28(4pt2):1421–1430.
- [54] Gonseth S, de Smith AJ, Roy R, et al. Genetic contribution to variation in DNA methylation at maternal smoking-sensitive loci in exposed neonates. *Epigenetics*. 2016 Sep;11(9):664–673.
- [55] Hannon E, Knox O, Sugden K, et al. Characterizing genetic and environmental influences on variable DNA methylation using monozygotic and dizygotic twins. *PLoS Genet*. 2018 Aug;14(8):e1007544.
- [56] Aryee MJ, Jaffe AE, Corrada-Bravo H, et al. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. *Bioinformatics*. 2014;30(10):1363–1369.
- [57] Maksimovic J, Gordon L, Oshlack A. SWAN: subset-quantile within array normalization for illumina infinium humanmethylation450 beadchips. *Genome Biol*. 2012 Jun;13(6):R44.
- [58] Du P, Zhang X, Huang -C-C, et al. Comparison of Beta-value and M-value methods for quantifying methylation levels by microarray analysis. *BMC Bioinformatics*. 2010;11:587.
- [59] Tandon A, Patterson N, Reich D. Ancestry informative marker panels for African Americans based on subsets of commercially available SNP arrays. *Genet Epidemiol*. 2011 Jan;35(1):80–83.
- [60] †. I. H. Consortium. The international HapMap project. *Nature*. 2003 Dec;426:789–796.
- [61] McKennan C, Nicolae D. Accounting for unobserved covariates with varying degrees of estimability in high-dimensional biological data. *Biometrika*. 2019 Sep;106(4):823–840. 0006–3444.

- [62] Storey JD. A direct approach to false discovery rates. *J R Stat Soc B*. 2001;63(3):479–498.
- [63] McKennan C, Nicolae D. Estimating and accounting for unobserved covariates in high-dimensional correlated data. *J Am Stat Assoc*. 2020 May;1–12.
- [64] Flutre T, Wen X, Pritchard J, et al. A statistical framework for joint eQTL analysis in multiple tissues. *PLoS Genet*. 2013 May;9(5):e1003486.
- [65] Stephens M. False discovery rates: a new deal. *Biostatistics*. 2017;18(2):275–294.