

Pruning and thresholding approach for methylation risk scores in multi-ancestry populations

Junyu Chen¹, Evan Gatev², Todd Everson^{3,1}, Karen N. Conneely⁴, Nastassja Koen^{5,6,7}, Michael P. Epstein⁴, Michael S. Kobor^{2,8,9}, Heather J. Zar^{10,11}, Dan J. Stein^{5,6,7}, Anke Huels^{1,3}

¹Department of Epidemiology, Rollins School of Public Health, Emory University, Atlanta, GA, USA

²Department of Medical Genetics, University of British Columbia, Vancouver, Canada

³Gangarosa Department of Environmental Health, Rollins School of Public Health, Emory University, Atlanta, Georgia, USA

⁴Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, USA

⁵Neuroscience Institute, University of Cape Town, Cape Town, South Africa

⁶Department of Psychiatry and Mental Health, University of Cape Town, Cape Town, South Africa

⁷South African Medical Research Council (SAMRC) Unit on Risk and Resilience in Mental Disorders, University of Cape Town, Cape Town, South Africa

⁸BC Children's Hospital Research Institute, Vancouver, Canada

⁹Centre for Molecular Medicine and Therapeutics, Vancouver, Canada

¹⁰Department of Pediatrics and Child Health, Red Cross War Memorial Children's Hospital, University of Cape Town, Cape Town, South Africa

¹¹South African Medical Research Council (SAMRC) Unit on Child and Adolescent Health, University of Cape Town, Cape Town, South Africa

Corresponding author:

Anke Hüls, PhD

Department of Epidemiology and Gangarosa Department of Environmental Health
Rollins School of Public Health, Emory University

1518 Clifton Road, Atlanta, GA 30322, USA

E-mail: anke.huels@emory.edu

Key Words: Epigenetic scores, Polygenic DNA methylation, Clumping and thresholding, Admixed population

Funding:

The Drakenstein Child Health Study was funded by the Bill & Melinda Gates Foundation (OPP 1017641), Discovery Foundation, South African Medical Research Council, National Research Foundation South Africa, CIDRI Clinical Fellowship and Wellcome Trust (204755/2/16/z). Additional support for the DNA methylation work was by the Eunice Kennedy Shriver National Institute of Child Health and Human Development of the National Institutes of Health (NICHD) under Award Number R21HD085849, and the Fogarty International Center (FIC). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. AH is supported the HERCULES Center (NIEHS P30ES019776). DJS and HJZ are supported by the South African Medical Research Council (SAMRC). The funders had no role in the study design, data collection and analysis, decision to publish, or preparation of manuscript.

Acknowledgments:

The authors thank the study and clinical staff at Paarl Hospital, Mbekweni and TC Newman clinics, as well as the CEO of Paarl Hospital, and the Western Cape Health Department for their support of the study. The authors thank the families and children who participated in this study.

Competing financial interests declaration:

All authors declare they have no actual or potential competing financial interest.

Abstract

Motivation: Recent efforts have focused on developing methylation risk scores (MRS), a weighted sum of the individual's DNAm values of pre-selected CpG sites.

Current MRS approaches only include genome-wide significant CpG sites and do not consider co-methylation. New methods that relax the p-value threshold to include more CpG sites and account for the inter-correlation of DNAm might improve the predictive performance of MRS.

Results: We paired informed co-methylation pruning with P-value thresholding to generate P+T MRS and evaluated its performance among multi-ancestry populations. Through simulation studies and real data analyses, we demonstrated that P+T MRS provides a substantial improvement over current thresholding methods for prediction of phenotypes. We demonstrated that European-derived summary statistics can be used to develop P+T MRS among other population such as African population. However, the prediction accuracy of P+T MRS may differ across multi-ancestry population due to environmental/cultural/social differences.

Availability and implementation: <https://github.com/jche453/Pruning-Thresholding-MRS>

Introduction

DNA methylation (DNAm), one of the most studied epigenetic mechanisms, regulates the mode of expression of DNA segments independent of alterations of their sequence by adding a methyl group at cytosine residues, hence contributing to variation in cellular phenotypes¹. With current advances in and reduction of cost of array-based profiling technologies, increasing numbers of large-scale epigenome-wide association studies (EWAS) have been conducted to study DNAm in association with complex human diseases as well as environmental and social factors^{2,3}. EWAS has thus far been successful in identifying dozens of cytosine guanine dinucleotides (CpGs) associated with various diseases and exposures, which could potentially be used for disease diagnosis and prediction, development of drug targets, and monitoring of drug response³⁻⁸. However, differential DNAm in individual CpGs often shows a weak prediction capacity and can only explain a small fraction of phenotype variance. Polyepigenetic approaches that aggregate information on differential DNAm from multiple CpGs might produce a more accurate biomarker for clinical usage^{9,10}.

The most popular polygenic approach for genotype data are polygenic risk scores (PRS), which are weighted sum of risk alleles of a pre-selected number of genetic variants¹¹. Recently, many efforts have focused on transferring PRS approaches to DNA methylation data to construct methylation risk scores (MRS), which are defined as weighted sums of the individuals' DNAm values of a pre-selected number of

CpGs¹⁰. However, there are many methodological challenges in constructing DNA methylation risk scores^{10,12,13}. One of the problems is that DNAm is influenced by ancestry, which is usually captured as self-reported race or ethnicity, and captures genetic ancestry (differences in the genome related to ancestry) as well as social determinants of health such as racism and discrimination, socioeconomic status, and environmental effects¹⁴. Thus, ideally, when external weights are used for the calculation of MRS, these weights should be assessed in a population with the same ancestry as the study samples. However, current epigenetic literature remains limited by the lack of diversity, with most focusing on European populations¹⁵, therefore making it difficult to identify appropriate weights for MRS for other populations. While it is well known that PRSs are not applicable across different ancestries^{16,17}, little is known about the performance of MRS across multi-ancestry populations.

Further, most of the current MRS approaches include subsets of individual CpGs, usually consisting of those that reached genome-wide significance in previously published epigenome-wide association studies (EWAS)¹⁰. Research in PRS showed that the optimal p-value threshold strongly depends on the data¹⁸, and including a larger proportion of variants could potentially capture more of the phenotype variation¹⁹. Moreover, as accounting for high linkage disequilibrium (LD) of single nucleotide polymorphisms (SNPs) improves the prediction performance of PRS²⁰, we hypothesize that accounting for DNA co-methylation, defined as proximal CpGs with

correlated DNAm across individuals²¹, could increase the prediction accuracy of MRS. One of the most popular PRS approaches to deal with single nucleotide polymorphisms (SNPs) in high linkage disequilibrium (LD) and to identify p-value thresholds with the best prediction accuracy is the pruning and threshold (P+T) method²². In the P+T approach, the correlation square (R^2) for SNPs within a close genetic distance is calculated and less significant SNPs that are correlated with an R^2 greater than a particular value (LD pruning)²³ are removed. Next, several p-value thresholds are tested to maximize the prediction accuracy of the derived PRS (p-value thresholding)^{23,24}. Theoretically, the P + T approach could be applied to generate MRS, however, there is no gold-standard on how to conduct pruning for DNAm data and the performance of such MRS across multi-ancestry populations remains unknown.

Here, we propose to use the Co-Methylation with genomic CpG Background (CoMeBack) method, a method that estimates DNA co-methylation, to account for correlations of DNAm at proximal CpG sites²¹, and pair it with p-value thresholding to construct P+T MRS. We conducted simulation studies based on data from an adult population consisting of three groups of different ancestries (Indian, White and Black, $n = 1,199$). Next, we applied the P+T approach to DNAm data from the Drakenstein Child Health Study ($n=270$)²⁵, a multi-ancestry birth cohort from South Africa, to evaluate the performance of MRS for maternal smoking status. Our simulation study and real data application demonstrate that the P+T approach improves the predictive

accuracy of MRS over the existing MRS methods. We also showed that MRS built upon the data from a population of one genetic ancestry could achieve high prediction performance among populations of other genetic ancestries, but the performance might differ in the presence of environmental/cultural/social differences associated with ancestry.

Materials and methods

P+T approach for MRS

P+T method refers to the calculation of MRS using informed co-methylation pruning (P) and P-value thresholding (T). First, summary statistics from an EWAS (typically include the participant ID, effect size, standard error and P-value of each CpG site) need to be estimated in an independent dataset (training dataset) to avoid overfitting, and then applied to generate MRS in a testing dataset (the samples used to evaluate the performance of MRS).

In our P+T method, co-methylation pruning is done by applying CoMeBack to DNAm data of the testing dataset or a reference panel²¹. Specifically, CoMeBack chains two adjacent array probes if the following requirements are met: 1. two probes are less than 2kb apart; 2. the reference human genome annotation shows a set of unmeasured genomic CpGs between them; 3. the density of unmeasured genomic CpGs between them is at least one CpG every 400bp. Chaining of adjacent array probes continues until an array probe does not meet the requirements, which will

form a unit where multiple CpGs are chained together. Correlations between DNAm levels will then be calculated for all array probes inside each unit. If all pairs of adjacent probes in a unit have a correlation square (R^2) greater than 0.3, such unit will be declared as a co-methylated region (CMR). Pruning is conducted by only keeping one CpG site per CMR in the dataset, the one with the lowest (most significant) P-value in the EWAS summary statistics.

Next, P-value thresholding step (T) is performed for the pruned set of CpG sites. Specifically, the P-value thresholding step (T) is performed by applying different P-value thresholds (e.g., P-value thresholds $\in [0.05, 0.005, 5 \times 10^{-4}, 5 \times 10^{-5}, \dots]$) and only including those CpG sites in the final MRS calculation that reached a P-value below those thresholds in the EWAS summary statistics.

Finally, for each P-value threshold, MRS are calculated as a weighted sum of DNAm β values (β value = methylated allele intensity / (unmethylated allele intensity + methylated allele intensity + 100), ranging from 0 representing unmethylation to 1 for complete methylation) of the selected CpGs, where the weights are the corresponding effect sizes for each CpG from the EWAS summary statistics. The squared correlation (R^2) between the phenotype of interest and MRS obtained using each P-value threshold is calculated to represent the prediction accuracy. The P-value threshold that produces MRS with the highest prediction accuracy in the testing data set is selected as the optimal P-value and the corresponding MRS is

used for downstream analysis. The pipeline for generating P+T MRS is written in an R script, which is available at GitHub (<https://github.com/jche453/Pruning-Thresholding-MRS.git>).

In our simulation studies and real data application, we compare the P+T MRS approach to the T approach, which refers to an approach in which the MRS is calculated by only using the thresholding approach described above, not accounting for correlations between included CpGs (no use of CoMeBack).

Simulation studies

To validate the performance of the proposed P+T MRS approach, we conducted simulation studies based on whole blood Illumina Infinium Human Methylation 450K BeadChip data from a discovery cohort composed of publicly available datasets. After the datasets were processed as previously described²¹, there were 1,199 adults (898 Indians, 136 Blacks and 165 Whites) and 386,362 CpGs left for MRS analysis. CoMeBack was applied to the DNA methylation β values of the 386,362 CpGs to obtain CMR. In each simulation, 10 of the 386,362 CpGs were randomly selected to be causal, k% (k = 30, 50, 70 or 100) of which are in a CMR with other CpGs. At most, one CpG would be causal in each CMR.

The causal CpGs were randomly assigned a “true” effect size from a uniform distribution as $w_i \sim U(-0.5, 0.5)$. We then simulated a phenotype for the j -th subject as follow:

$$Y_j = \sum_{i=1}^{10} w_i m_{ij} + \varepsilon_j, \varepsilon_j \sim N(0, \delta^2),$$

where m_{ij} is the DNAm β value of causal CpG site i of the j -th subject, and ε_j is an error term that follows a normal distribution. Different δ^2 were set to ensure that the targeted variance of phenotype explained by DNAm alone equals 10%, 30% 50% or 80%.

We also simulated a second phenotype Y_j^* for the j -th subject, which was directly affected by ancestry using European ancestry as reference:

$$Y_j^* = \sum_{i=1}^{10} w_i m_{ij} + a * (\text{if Indian}) + b * (\text{if Black}) + \varepsilon_j^*, \varepsilon_j^* \sim N(0, \delta^{*2})$$

Effect of ancestry in our simulations is simulated as the effect of genetic ancestry assuming there were no complex social determinants involved in the causal pathway. Different δ^{*2} were used so that the variance of phenotype that was explained by DNAm and ancestry together equals 20%, 50% or 80. For our simulations, effect a was set to 0.1 and b to 0.2. In each simulation, both simulated phenotypes Y_j and Y_j^* share the same epigenetic liability ($\sum_{i=1}^{10} w_i m_{ij}$).

In each simulation, for fair comparison, 762 Indians were randomly chosen as the training dataset so that there were at least 136 people left for each race group in the testing dataset. Associations between CpGs and each of the two simulated

phenotypes were assessed by robust linear regression model using limma R package²⁶ in the training dataset. We calculated top 10 principal components (PCs) from DNAm of 386,362 CpGs²⁷ and used EpiDISH to estimate cell type proportions of each CpGs²⁸. We observed that in our simulation dataset, top 10 PCs are highly correlated with cell type proportions (Supplement figure 1A), and using either summary statistics adjustment for top 10 PCs or summary statistics adjusted for cell type proportions would lead to almost identical prediction performance of MRS (Supplement figure 1B). Thus, to account for population stratification and cell type difference, we adjusted for the top 10 PCs in our main analyses. The summary statistics (effect size and P-values) obtained from association tests in the training data were saved and later used to construct MRS in the testing dataset. We repeated 1000 simulations per scenario to evaluate the prediction accuracy (R^2), power and type 1 error rate of the P+T MRS.

We evaluated the performance of the MRS not only in scenarios of A) same ancestry in training and test data, but also B) across different ancestry groups (training data: Indian, test data: European or African) and C) in multi-ancestry populations (training data: Indian, test data: Indian, European and African). For scenario C), we evaluated two analysis strategies: 1. Joint-analysis: perform MRS analyses in the whole testing dataset where subjects from all racial groups were merged; 2. Standardization: scale MRS to have a standard normal distribution within each racial group before merging all subjects for analyses.

Application study of smoking MRS

To evaluate the performance of the P+T approach in a real data setting, we applied the P+T approach to calculate a MRS for maternal smoking status during pregnancy using cord blood DNAm data from newborns in the South African Drakenstein Child Health Study (DCHS), a multi-ancestry longitudinal study investigating determinants of early child development²⁹. There were 145 Black African infants and 115 Mixed ancestry infants in the DCHS. A detailed description of the enrollment process, inclusion criteria, variables measurement and ethical approval of the study have been previously published^{29,30}.

Cotinine levels were measured in urine provided by mothers within four weeks of enrollment and classified as <10 ng/ml (non-smoker), 10–499 ng/ml (passive smoker), or ≥500 ng/ml (active smoker)²⁵. We also considered a dichotomized smoking variable indicating if the pregnant women were active smoker versus non-active smoker (passive smoker + non-smoker) in evaluation of the prediction performance of MRS. Cord blood was collected at time of delivery and used to measure DNA methylation by either MethylationEPIC BeadChips (EPIC, n=145) or the Illumina Infinium HumanMethylation450 BeadChips (450K, n=103)^{29,30}, followed by quality control and normalization to calculate β values (details have been published elsewhere)³¹.

Summary statistics for the calculation of MRS were obtained from a study that meta-analyzed the associations between newborn blood DNA methylation and sustained maternal smoking during pregnancy among 5,648 mother-child pairs as part of the Pregnancy and Childhood Epigenetics (PACE) Consortium (Table 1)³². The participants of all cohorts used in the meta-analysis except one were of European ancestry.

In addition, we compared P+T MRS to three previously published MRS for maternal smoking during pregnancy (Reese MRS, Richmond 19 MRS, Richmond 568 MRS; Table 1). Reese MRS model was trained among 1,068 newborns of European ancestry in the Norwegian Mother and Child Cohort Study, while Richmond 568 MRS and Richmond 19 MRS was trained in multi-ancestry newborns (N=6,685) and children around 6.8 years old (N=3,187) in PACE Consortium respectively. The training population for Reese MRS and Richmond 19 MRS overlapped with the training population for summary statistics used in P+T MRS in our study. Reese et al. used a LASSO regression to select CpGs for Reese MRS, which is a weighted sum of DNAm β values of 28 CpGs with weights estimated from the LASSO regression³³. Richmond 19 MRS is a weighted sum of DNAm β values of 19 CpGs that were significantly associated with prenatal smoking in an EWAS conducted in peripheral blood from children of averaged 6.8 years age (Richmond 19 MRS)^{34,35}. In the same study, Richmond 568 MRS was proposed based on 568 CpGs that were significantly associated with prenatal smoking in cord blood^{34,35}.

Results

Simulation results

We compared the prediction performance of P+T MRS to the T method among 136 Indians in the test data across different simulation scenarios (Figure 1). Figure 1A shows that P+T MRS that account for co-methylation between CpGs have stable prediction performance when proportion of causal CpGs located in a CMR (k%) varies. In contrast, the T method had a slightly lower prediction performance with a wider range compared to P+T MRS when k% is small and the prediction accuracy decreased with increasing k%. While the P+T MRS outperformed T method when V_{DNAm}^2 is 80% (Figure 1A), the difference between P+T MRS and the T method decreases as the V_{DNAm}^2 decreases. This is likely because as variance explained by DNAm decreases, there is less power for association testing, and it becomes increasingly difficult to distinguish real signals from statistical noise while generating the summary statistics.

Next, we assessed the performance of P+T and T method across different ancestries and among multi-ancestry populations. Both P+T MRS and T method achieved a high power (> 95%) and a low type 1 error rate (~ 5%) within each ancestry in most scenarios for both phenotypes except when the phenotype variance explained by DNA methylation is 10% or 30% (Supplement table 1-4). When the phenotypes are not associated with ancestry (Figure 2A), the three MRS analyses strategies

(stratification, joint analysis and standardization) lead to nearly identical results.

However, when the phenotypes are ancestry-dependent, both joint analysis and standardization of MRS showed very poor prediction of the phenotypes (Figure 2B).

Findings were similar when the phenotype variance explained by DNA methylation was reduced from 80% to 10%, 30% or 50% (Supplement Figure 2). On the contrary, when the phenotype variance explained by DNA methylation was low, joint analysis and standardization of MRS increased the power due to increased sample sizes in comparison to the analyses of the individual subgroups (Supplement table 3).

MRS of maternal smoking status

Figure 3 shows the prediction performance of MRS for maternal smoking status among DCHS newborns. As the p-value threshold decreases, the prediction accuracy of the resulting MRS increases before reaching a plateau, demonstrating the importance of P-value thresholding in MRS to control for noise. Among mixed ancestry newborns, P+T MRS of smoking status achieved a prediction accuracy of 29.5% using P-value threshold of 5×10^{-21} , while the best T method MRS have a lower prediction accuracy (24.5%) using P-value threshold 5×10^{-10} , confirming the benefits of pruning in MRS calculation (Figure 3A). Both P+T MRS and T method MRS had lower prediction performance for maternal smoking among Black African infants (10.9% and 8.0% respectively) (Figure 3B), which is likely due to the low prevalence of active smoking among mothers of Black African infants in DCHS (13%) compared to mothers of mixed ancestry infants (49%) (Supplement figure 3).

Joint analysis of P+T MRS showed a prediction accuracy of 20.4%, which is between the prediction accuracy of P+T MRS among Black African infants and mixed ancestry infants (Figure 3C). Standardization approach did not improve the performance of MRS (Figure 3D).

We next compared the prediction accuracy and distribution of P+T MRS to other established MRS for maternal smoking during pregnancy and newborn DNAm (Figure 4). P+T MRS showed a similar prediction accuracy as Reese MRS among both Black and Mixed ancestry infants, and they outperformed all other smoking MRS (Figure 4A). The distributions of P+T MRS in Black African infants and mixed ancestry infants were similar within each category of maternal smoking status (Figure 4B), whereas the other MRS had different distributions in Black infants and Mixed ancestry infants within each category of smoking, suggesting that P+T MRS might be more comparable across different ancestries than the 3 established MRS approaches (Figure 4C, Supplement figure 3). Furthermore, the P+T approach is computationally efficient and does not require access to individual DNAm data from the training data since it makes use of published EWAS summary statistics.

Discussion

Based on the well-established P+T framework in PRS, we developed P+T MRS, which aggregates EWAS signals and could potentially be used as a biomarker in association studies where single CpGs do not achieve significance^{4,36,37}. The

proposed P+T MRS approach uses CoMeBack for co-methylation pruning and evaluates multiple P-value thresholds to maximize prediction performance. Such MRS could potentially serve as a powerful dimension reduction approach for mediation and multi-omics integration analyses^{4,36-39} as well as biomarkers of individual disease risk in a clinical setting⁴⁰⁻⁴².

Overall, our simulation studies demonstrated good performance of P+T MRS for predicting phenotypes of interest with good statistical power and well-controlled type 1 error. We demonstrated that the prediction accuracy of MRS reflects the variance of phenotype that is explained by DNAm. By accounting for inter-correlation between CpGs, P+T MRS showed a better performance than the standard T method, especially when a large number of causal CpGs are located in CMRs. In the real data application, we further confirmed the improvement of P+T MRS, with improved prediction of maternal smoking status compared to the T method. However, P+T MRS could still have poor prediction performance if the external EWAS is underpowered or subject to bias.

In the prediction of maternal smoking status, P+T MRS showed comparable performance to Reese MRS, which was derived using the LASSO method³³. When predictors are highly correlated, LASSO typically selects one of the correlated predictors and shrinks the effect size of the rest to zero, which might produce similar results to our pruning procedure in developing MRS. One of the advantages of P+T

MRS is that it is based on EWAS summary statistics which are often publicly available, making it a valuable approach, as it is often difficult to obtain individual DNAm data from an external cohort. In contrast, to construct MRS like Reese MRS, individual DNAm data are usually required to perform a LASSO regression, and these are often not accessible. Recently, novel penalized regressions have been proposed to generate PRS with only GWAS summary statistics and publicly available reference data⁴³, but their applications to EWAS summary statistics for MRS have not been investigated. To develop MRS for different exposures and outcomes, we urge EWAS studies to make their genome-wide summary statistics publicly available.

In our simulation studies, weights obtained from Indian training samples were applied to generate P+T MRS, thus MRS among Indian testing samples were assumed to have the best prediction of the simulated phenotypes. However, MRS among Whites and Blacks also achieved a prediction accuracy as high as among Indians for both simulated phenotypes suggesting that genetic ancestry does not contribute to difference in prediction abilities of MRS across multi-ancestry population. This is likely because we assumed all ancestries share the same causal CpGs and effect sizes. However, in the real world, this assumption could possibly be violated for many phenotypes. Furthermore, unlike ancestry in our simulation studies, ancestry in the real world is complex. The meaning of ancestry could be different in different regions/nations, and “effect of ancestry” involves the joint effects

of ancestry-associated social determinants of health and environmental effects, and cultural context⁴⁴. Ancestry, along with environment and social differences associated with it, could affect both MRS and phenotypes in numerous causal pathways and potentially modify the effect of MRS on the phenotypes. Thus, even if all ancestries indeed share the same causal CpGs and effect sizes, it might still not be sufficient to disentangle the relationship between ancestry, DNAm and phenotype of interest. This may greatly impact the transferability of MRS across different ancestries, which could be the reason why we observed an inconsistency of performance of P+T MRS in terms of their distributions and predictions across multi-ancestry population in the real data analyses. In practice, we recommend that researchers conduct MRS analyses stratified by ancestry first and evaluate the effect of ancestry on MRS analyses before pooling participants together for a joint analysis.

In our real data application, summary statistics for smoking were obtained from a cohort with mainly people of European ancestry³². MRS of smoking among mixed ancestry infants achieved a prediction accuracy of nearly 30%. However, the prediction accuracy of P+T MRS among Black African infants was only 10.9%. We suspect that the difference was largely due to the prevalence of active smoking among mothers of Black African infants being lower than those of mixed ancestry infants (13% vs 49%), which is similar to how the prevalence of outcome affects the predictive ability of PRS⁴⁵.

To the best of our knowledge, this is the first study to evaluate P+T approach in the construction of MRS among multi-ancestry populations. However, there are several potential limitations that warrant mention. First, the sample size of both simulation studies and real data analyses was relatively small, thus our results might not fully capture the strengths and limitations of P+T MRS. Second, lack of different ancestry-specific summary statistics made it impossible to compare the use of external weights from population of different ancestries (e.g. European ancestry vs other ancestries). Third, the prevalence of active maternal smoking is different in different ancestries and has influenced the performance of P+T MRS. As a result, real data analysis of smoking MRS could not provide firm evidence about the transferability of MRS between Black African and mixed ancestry infants. Fourth, we mainly focused on the prediction performance of P+T MRS. Further studies are needed to assess the performance of P+T MRS in mediation analysis.

In conclusion, P+T MRS provides a substantial improvement for prediction of phenotype of interest, either exposures or diseases, over T method that does not account for co-methylation between CpGs. In contrast to PRS, using existing summary statistics that were derived from European populations can be used to calculate MRS in other ancestries, thus reducing the ancestry/ethnicity disparity in medical research. However, caution is needed in the analyses and interpretation of MRS results across multi-ancestry populations. More investigations of MRS are urged to further improve their prediction accuracy and translational values, also in

combination with other clinical and non-clinical variables, especially among multi-ancestry population. With the current increase of large consortia-led EWAS for different exposures and health outcomes (e.g., the PACE consortium), we believe the predictive performance of MRS will continue to increase, and the P+T method has the potential to be widely used for risk prediction and mediation analyses.

Reference

1. Berger SL, Kouzarides T, Shiekhhattar R, Shilatifard A. An operational definition of epigenetics. *Genes Dev.* Apr 1 2009;23(7):781-3. doi:10.1101/gad.1787609
2. Rakyan VK, Down TA, Balding DJ, Beck S. Epigenome-Wide Association Studies for common human diseases. *Nature reviews Genetics.* 07/12 2011;12(8):529-541. doi:10.1038/nrg3000
3. Wei S, Tao J, Xu J, et al. Ten Years of EWAS. *Adv Sci (Weinh).* Aug 11 2021:e2100727. doi:10.1002/advs.202100727
4. Wahl S, Drong A, Lehne B, et al. Epigenome-wide association study of body mass index, and the adverse outcomes of adiposity. *Nature.* Jan 5 2017;541(7635):81-86. doi:10.1038/nature20784
5. McCartney DL, Hillary RF, Stevenson AJ, et al. Epigenetic prediction of complex traits and death. *Genome biology.* 2018;19(1):136-136. doi:10.1186/s13059-018-1514-1
6. Heiss JA, Brenner H. Epigenome-wide discovery and evaluation of leukocyte DNA methylation markers for the detection of colorectal cancer in a screening setting. *Clinical epigenetics.* 2017;9:24-24. doi:10.1186/s13148-017-0322-x
7. Gelato KA, Shaikhibrahim Z, Ocker M, Haendler B. Targeting epigenetic regulators for cancer therapy: modulation of bromodomain proteins, methyltransferases, demethylases, and microRNAs. *Expert Opin Ther Targets.* Jul 2016;20(7):783-99. doi:10.1517/14728222.2016.1134490
8. Krushkal J, Silvers T, Reinhold WC, et al. Epigenome-wide DNA methylation analysis of small cell lung cancer cell lines suggests potential chemotherapy targets. *Clinical epigenetics.* 2020;12(1):93-93. doi:10.1186/s13148-020-00876-8
9. Guan Z, Yu H, Cuk K, Zhang Y, Brenner H. Whole-Blood DNA Methylation Markers in Early Detection of Breast Cancer: A Systematic Literature Review. *Cancer Epidemiol Biomarkers Prev.* Mar 2019;28(3):496-505. doi:10.1158/1055-9965.Epi-18-0378
10. Hüls A, Czamara D. Methodological challenges in constructing DNA methylation risk scores. *Epigenetics.* Jan-Feb 2020;15(1-2):1-11. doi:10.1080/15592294.2019.1644879
11. Wray NR, Lee SH, Mehta D, Vinkhuyzen AA, Dudbridge F, Middeldorp CM. Research review: Polygenic methods and their application to psychiatric traits. *J Child Psychol Psychiatry.* Oct 2014;55(10):1068-87. doi:10.1111/jcpp.12295
12. Martin EM, Fry RC. Environmental Influences on the Epigenome: Exposure-Associated DNA Methylation in Human Populations. *Annual Review of Public Health.* 2018/04/01 2018;39(1):309-333. doi:10.1146/annurev-publhealth-040617-014629
13. Notterman DA, Mitchell C. Epigenetics and Understanding the Impact of Social Determinants of Health. *Pediatr Clin North Am.* Oct 2015;62(5):1227-40. doi:10.1016/j.pcl.2015.05.012
14. Borrell LN, Elhawary JR, Fuentes-Afflick E, et al. Race and Genetic Ancestry in Medicine — A Time for Reckoning with Racism. *New England Journal of Medicine.* 2021/02/04 2021;384(5):474-480. doi:10.1056/NEJMms2029562
15. Cronjé HT, Elliott HR, Nienaber-Rousseau C, Pieters M. Replication and expansion of epigenome-wide association literature in a black South African population. *Clinical Epigenetics.* 2020/01/07 2020;12(1):6. doi:10.1186/s13148-019-0805-z

16. Khera AV, Chaffin M, Zekavat SM, et al. Whole-Genome Sequencing to Characterize Monogenic and Polygenic Contributions in Patients Hospitalized With Early-Onset Myocardial Infarction. *Circulation*. Mar 26 2019;139(13):1593-1602. doi:10.1161/circulationaha.118.035658
17. Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nature Genetics*. 2019/04/01 2019;51(4):584-591. doi:10.1038/s41588-019-0379-x
18. Goldstein BA, Yang L, Salfati E, Assimes TL. Contemporary Considerations for Constructing a Genetic Risk Score: An Empirical Approach. *Genet Epidemiol*. Sep 2015;39(6):439-45. doi:10.1002/gepi.21912
19. Euesden J, Lewis CM, O'Reilly PF. PRSice: Polygenic Risk Score software. *Bioinformatics (Oxford, England)*. 2015;31(9):1466-1468. doi:10.1093/bioinformatics/btu848
20. Choi SW, Mak TS, O'Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. *Nat Protoc*. Sep 2020;15(9):2759-2772. doi:10.1038/s41596-020-0353-1
21. Gatev E, Gladish N, Mostafavi S, Kobor MS. CoMeBack: DNA methylation array data analysis for co-methylated regions. *Bioinformatics*. May 1 2020;36(9):2675-2683. doi:10.1093/bioinformatics/btaa049
22. Wu J, Pfeiffer RM, Gail MH. Strategies for developing prediction models from genome-wide association studies. *Genet Epidemiol*. Dec 2013;37(8):768-77. doi:10.1002/gepi.21762
23. Privé F, Vilhjálmsson BJ, Aschard H, Blum MGB. Making the Most of Clumping and Thresholding for Polygenic Scores. *The American Journal of Human Genetics*. 2019/12/05/ 2019;105(6):1213-1221. doi:<https://doi.org/10.1016/j.ajhg.2019.11.001>
24. Choi SW, O'Reilly PF. PRSice-2: Polygenic Risk Score software for biobank-scale data. *GigaScience*. 2019;8(7)doi:10.1093/gigascience/giz082
25. Vanker A, Barnett W, Workman L, et al. Early-life exposure to indoor air pollution or tobacco smoke and lower respiratory tract illness and wheezing in African infants: a longitudinal birth cohort study. *Lancet Planet Health*. Nov 2017;1(8):e328-e336. doi:10.1016/s2542-5196(17)30134-1
26. Ritchie ME, Phipson B, Wu D, et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res*. Apr 20 2015;43(7):e47. doi:10.1093/nar/gkv007
27. Barfield RT, Almli LM, Kilaru V, et al. Accounting for population stratification in DNA methylation studies. *Genet Epidemiol*. Apr 2014;38(3):231-41. doi:10.1002/gepi.21789
28. Teschendorff AE, Breeze CE, Zheng SC, Beck S. A comparison of reference-based algorithms for correcting cell-type heterogeneity in Epigenome-Wide Association Studies. *BMC Bioinformatics*. Feb 13 2017;18(1):105. doi:10.1186/s12859-017-1511-5
29. Zar HJ, Barnett W, Myer L, Stein DJ, Nicol MP. Investigating the early-life determinants of illness in Africa: the Drakenstein Child Health Study. *Thorax*. Jun 2015;70(6):592-4. doi:10.1136/thoraxjnl-2014-206242
30. Stein DJ, Koen N, Donald KA, et al. Investigating the psychosocial determinants of child health in Africa: The Drakenstein Child Health Study. *J Neurosci Methods*. Aug 30 2015;252:27-35. doi:10.1016/j.jneumeth.2015.03.016

31. Hüls A, Wedderburn CJ, Groenewold NA, et al. Newborn differential DNA methylation and subcortical brain volumes as early signs of severe neurodevelopmental delay in a South African Birth Cohort Study. *World J Biol Psychiatry*. Jan 12 2022;1-12. doi:10.1080/15622975.2021.2016955
32. Sikdar S, Joehanes R, Joubert BR, et al. Comparison of smoking-related DNA methylation between newborns from prenatal exposure and adults from personal smoking. *Epigenomics*. Oct 2019;11(13):1487-1500. doi:10.2217/epi-2019-0066
33. Reese SE, Zhao S, Wu MC, et al. DNA Methylation Score as a Biomarker in Newborns for Sustained Maternal Smoking during Pregnancy. *Environ Health Perspect*. Apr 2017;125(4):760-766. doi:10.1289/ehp333
34. Richmond RC, Suderman M, Langdon R, Relton CL, Davey Smith G. DNA methylation as a marker for prenatal smoke exposure in adults. *International Journal of Epidemiology*. 2018;47(4):1120-1130. doi:10.1093/ije/dyy091
35. Joubert Bonnie R, Felix Janine F, Yousefi P, et al. DNA Methylation in Newborns and Maternal Smoking in Pregnancy: Genome-wide Consortium Meta-analysis. *The American Journal of Human Genetics*. 2016/04/07/ 2016;98(4):680-696. doi:<https://doi.org/10.1016/j.ajhg.2016.02.019>
36. Yu H, Raut JR, Schöttker B, Holleczeck B, Zhang Y, Brenner H. Individual and joint contributions of genetic and methylation risk scores for enhancing lung cancer risk stratification: data from a population-based cohort in Germany. *Clin Epigenetics*. Jun 18 2020;12(1):89. doi:10.1186/s13148-020-00872-y
37. Westerman K, Fernández-Sanlés A, Patil P, et al. Epigenomic Assessment of Cardiovascular Disease Risk and Interactions With Traditional Risk Metrics. *J Am Heart Assoc*. Apr 21 2020;9(8):e015299. doi:10.1161/jaha.119.015299
38. Guan Z, Raut JR, Weigl K, et al. Individual and joint performance of DNA methylation profiles, genetic risk score and environmental risk scores for predicting breast cancer risk. *Mol Oncol*. 2020;14(1):42-53. doi:10.1002/1878-0261.12594
39. Grant CD, Jafari N, Hou L, et al. A longitudinal study of DNA methylation as a potential mediator of age-related diabetes risk. *Geroscience*. Dec 2017;39(5-6):475-489. doi:10.1007/s11357-017-0001-z
40. García-Calzón S, Perfilyev A, Martinell M, et al. Epigenetic markers associated with metformin response and intolerance in drug-naïve patients with type 2 diabetes. *Sci Transl Med*. Sep 16 2020;12(561)doi:10.1126/scitranslmed.aaz1803
41. Deng Y, Wan H, Tian J, et al. CpG-methylation-based risk score predicts progression in colorectal cancer. *Epigenomics*. Apr 2020;12(7):605-615. doi:10.2217/epi-2019-0300
42. Kilanowski A, Chen J, Everson T, et al. Methylation risk scores for childhood aeroallergen sensitization: Results from the {LISA} birth cohort. *Authorea*. Nov 2021;doi:10.22541/au.163620398.85835627/v1
43. Pattee J, Pan W. Penalized regression and model selection methods for polygenic scores on summary statistics. *PLoS Comput Biol*. Oct 2020;16(10):e1008271. doi:10.1371/journal.pcbi.1008271
44. VanderWeele TJ, Robinson WR. On the causal interpretation of race in regressions adjusting for confounding and mediating variables. *Epidemiology*. Jul 2014;25(4):473-84. doi:10.1097/ede.000000000000105

45. Gibson G. On the utilization of polygenic risk scores for therapeutic targeting. *PLoS Genet.* Apr 2019;15(4):e1008060. doi:10.1371/journal.pgen.1008060

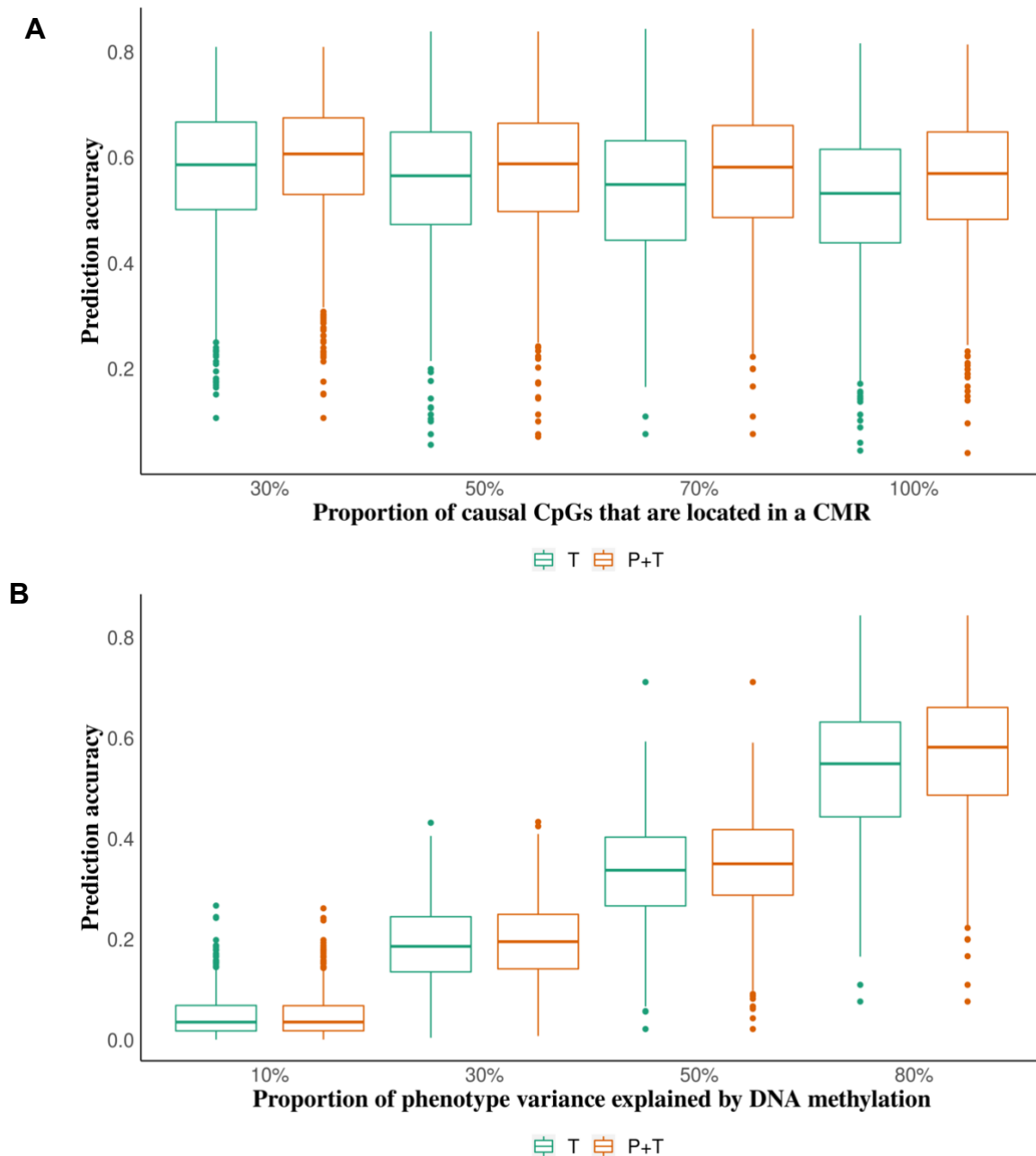


Figure 1. Simulation study. Prediction accuracy of P+T and T method in dependence of (A) the proportion of causal CpG sites in CMRs and (B) proportion of phenotype variance explained by DNA methylation, among Indian participants. For each simulation, the discovery cohort was repeatedly and randomly split into a training set comprising 762 Indians and a testing set comprising 136 people of the same ancestry. Phenotypes were simulated without an influence of ancestry. Results are shown for (A) different proportions of causal CpGs located in CMR (30%, 50%, 70%, 100%) and (B) different proportions of phenotype variance explained by DNA methylation (10%, 30%, 50%, 80%). Each box represents the distribution of prediction accuracy across 1000 simulations, where the central mark is the median and the edges of the box are the 25th and 75th percentiles.

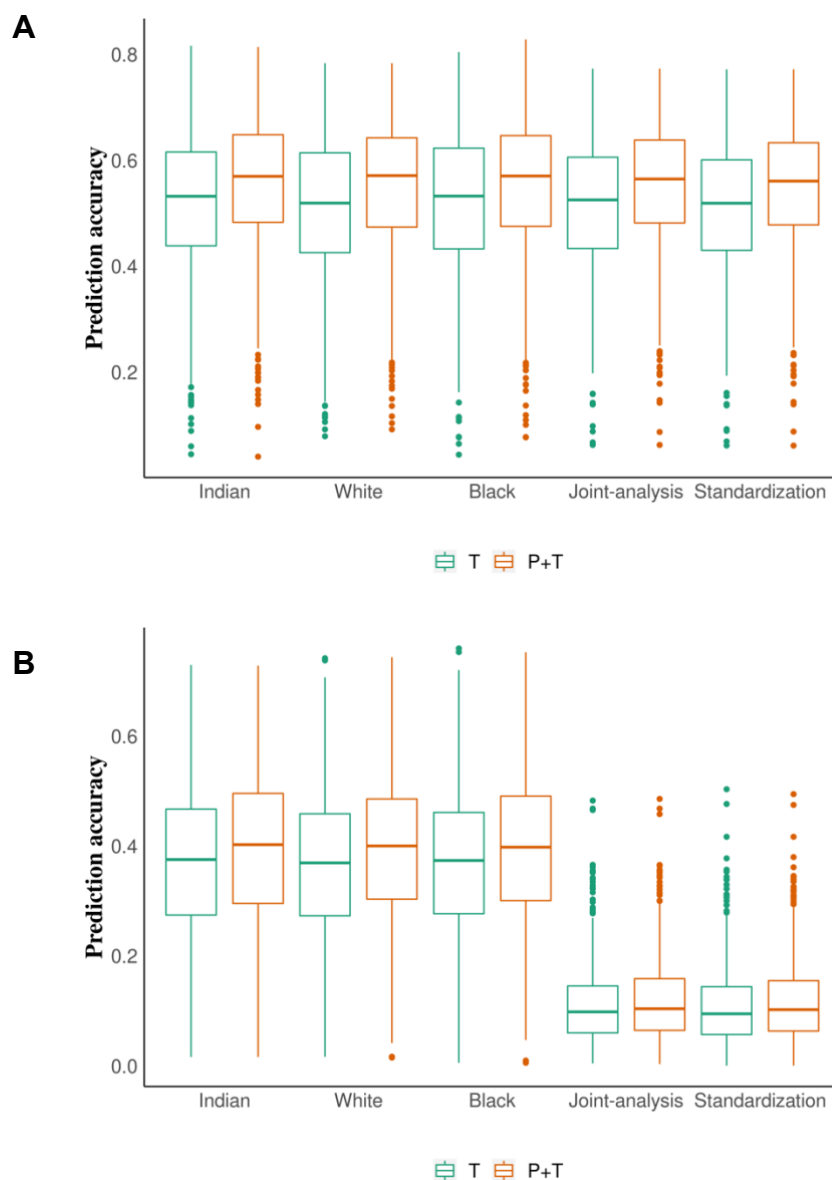


Figure 2. Simulation study. Prediction accuracy of P+T and T approach across different racial groups and among multi-ancestry populations. For each simulation, the discovery cohort was repeatedly and randomly split into a training set comprising 762 Indians and a testing set comprising 136 people of each ancestry group. The proportion of causal CpGs located in CMR is 70% and the proportion of phenotype variance explained by DNA methylation (and ancestry) is 80%. Results are shown for the prediction of simulated phenotypes (**2A**) without an influence of ancestry and (**2B**) influenced by ancestry. Joint-analysis refers to MRS analyses of all participants pooled from all ancestry groups and standardization refers to standardizing MRS within each ancestry group and then merging all participants before analyses. Each box represents the distribution of prediction accuracy across 1000 simulations, where the central mark is the median and the edges of the box are the 25th and 75th percentiles.

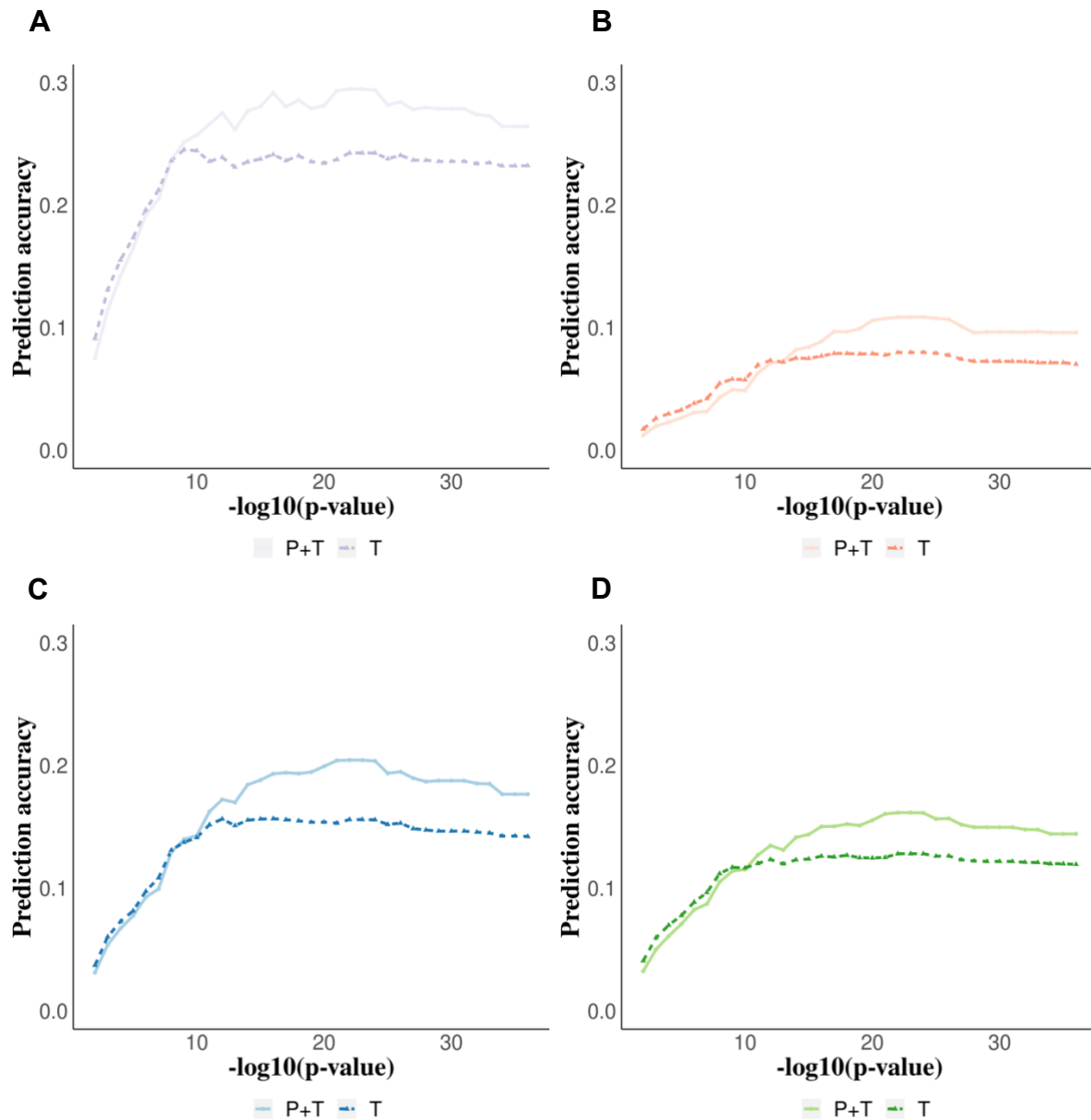


Figure 3. Real data application. MRS for the prediction of maternal smoking during pregnancy using cord blood DNA methylation data from newborns in the South African Drakenstein Child Health Study (DCHS). Prediction accuracy of maternal smoking status is shown stratified for **A. Mixed infants. **B.** Black infants. **C.** joint-analysis (all subjects pooled from all ancestries) **D.** Standardization (standardizing MRS within each ancestry and merging all subjects before analyses).**

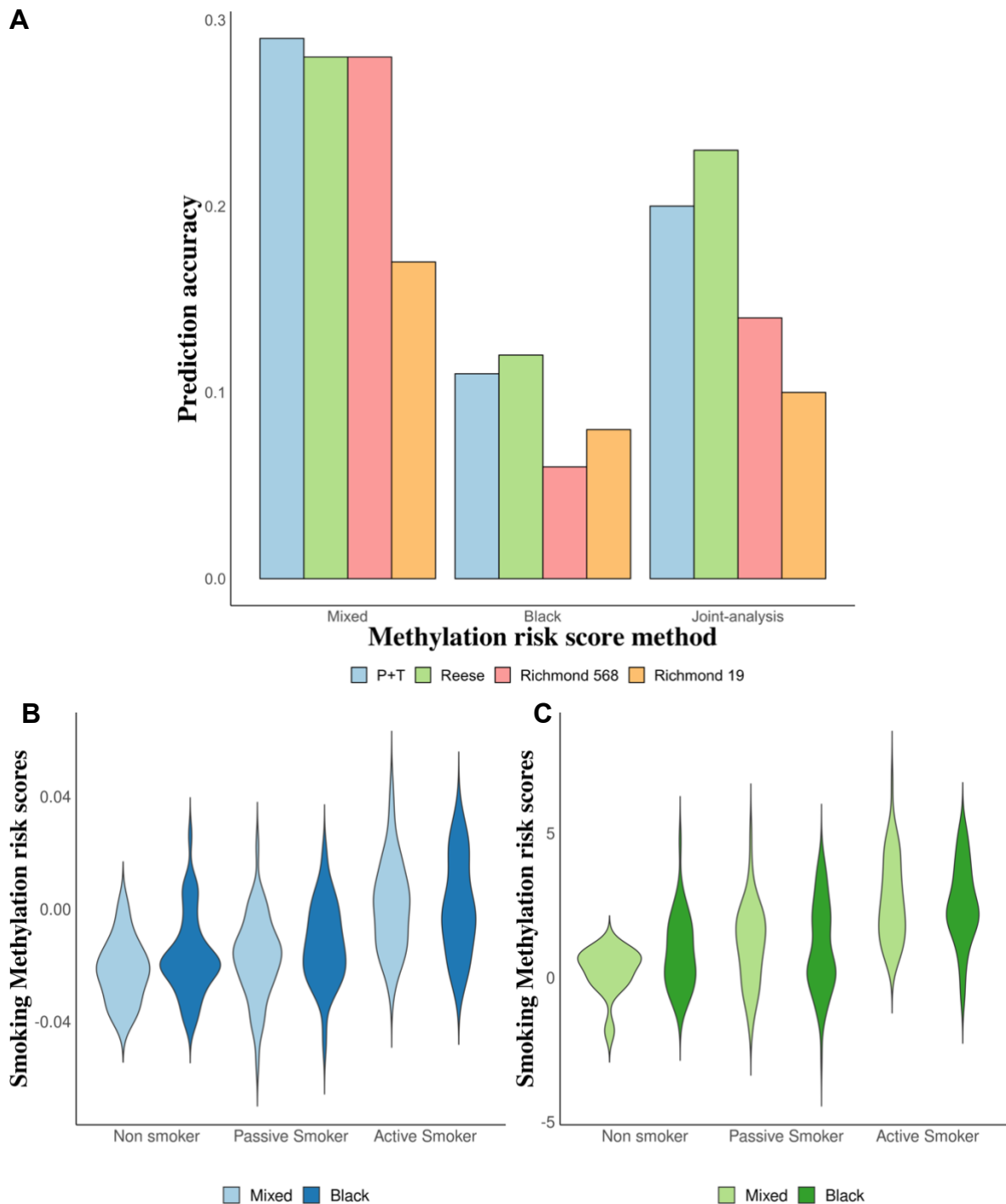


Figure 4. Real data application. Comparison of P+T method to 3 other published MRS for predicting maternal smoking status in the South African Drakenstein Child Health Study (DCHS). A Prediction performance of all 5 MRS methods for Mixed infants, Black infants and pooled samples (joint-analysis). Distribution of **(B)** P+T MRS and **(C)** Reese MRS (best performing MRS in **(A)**) among non-smokers, passive smokers and active smokers stratified by ancestry.

Table 1. Overview of included EWAS, their phenotypes, training sample and methods.

MRS	Training dataset publication	Training Population	Phenotype	MRS publication	Method	No. of CpG sites
P+T MRS	Sikdar et al. 2019 ¹	Multi-ethnic newborns (mainly White, N=5,648)	Most cohorts ascertained sustained smoking during pregnancy by questionnaires; two cohorts incorporated cotinine-based smoking measure	-	P+T	22
Reese MRS	Reese et al. 2017 ²	White newborns (N=1,068)*	Sustained smoking during pregnancy obtained from combined information of cotinine-based and self-report based classification	Reese et al. 2017 ²	Logistic LASSO regression	28
Richmond 568 MRS	Joubert et al. 2016 ³	Multi-ethnic newborns (N=6,685)*	Maternal smoking during pregnancy via questionnaires	Richmond et al. 2018 ⁴	Robust linear regression; Bonferroni corrected P-value < 0.05	568
Richmond 19 MRS	Joubert et al. 2016 ³	Multi-ethnic older children (average age = 6.8 years) (N=3,187)	Maternal smoking during pregnancy via questionnaires	Richmond et al. 2018	Robust linear regression; Bonferroni corrected P-value < 0.05	19

* These training populations overlapped with training population for summary statistics used for P+T MRS.

Reference

1. Sikdar S, Joehanes R, Joubert BR, et al. Comparison of smoking-related DNA methylation between newborns from prenatal exposure and adults from personal smoking. *Epigenomics*. Oct 2019;11(13):1487-1500. doi:10.2217/epi-2019-0066
2. Reese SE, Zhao S, Wu MC, et al. DNA Methylation Score as a Biomarker in Newborns for Sustained Maternal Smoking during Pregnancy. *Environ Health Perspect*. Apr 2017;125(4):760-766. doi:10.1289/ehp333
3. Joubert Bonnie R, Felix Janine F, Yousefi P, et al. DNA Methylation in Newborns and Maternal Smoking in Pregnancy: Genome-wide Consortium Meta-analysis. *The American Journal of Human Genetics*. 2016/04/07/ 2016;98(4):680-696. doi:<https://doi.org/10.1016/j.ajhg.2016.02.019>
4. Richmond RC, Suderman M, Langdon R, Relton CL, Davey Smith G. DNA methylation as a marker for prenatal smoke exposure in adults. *International Journal of Epidemiology*. 2018;47(4):1120-1130. doi:10.1093/ije/dyy091