

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^{8,9,10}, Heather J. Zar^{11,12}, Dan J. Stein^{5,6,7}, Anke Huels^{1,3}

1. Department of Epidemiology, Rollins School of Public Health, Emory University, Atlanta, GA, USA
2. Institute of Molecular Biology "Acad. Roumen Tsanev", Sofia, Bulgaria
3. Gangarosa Department of Environmental Health, Rollins School of Public Health, Emory University, Atlanta, Georgia, USA
4. Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, USA
5. Neuroscience Institute, University of Cape Town, Cape Town, South Africa
6. Department of Psychiatry and Mental Health, University of Cape Town, Cape Town, South Africa
7. South African Medical Research Council (SAMRC) Unit on Risk and Resilience in Mental Disorders, University of Cape Town, Cape Town, South Africa
8. Department of Medical Genetics, University of British Columbia, Vancouver, Canada
9. BC Children's Hospital Research Institute, Vancouver, Canada
10. Centre for Molecular Medicine and Therapeutics, Vancouver, Canada
11. Department of Pediatrics and Child Health, Red Cross War Memorial Children's Hospital, University of Cape Town, Cape Town, South Africa
12. South 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

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.

Disclosure statement:

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

Data Availability Statement:

The data for simulation studies were derived from the following dataset (GSE55763, GSE84727, GSE80417, GSE111629 and GSE72680) from NCBI Gene Expression Omnibus (<https://www.ncbi.nlm.nih.gov/geo/>). The data for real data application are available from the corresponding author, [AH], upon reasonable request.

Abstract

Recent efforts have focused on developing methylation risk scores (MRS), a weighted sum of the individual's DNAm values of pre-selected CpG sites. Most of the current MRS approaches that utilize Epigenome-wide association studies (EWAS) summary statistics 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. We paired informed co-methylation pruning with P-value thresholding to generate pruning and thresholding (P+T) MRS and evaluated its performance among multi-ancestry populations. Through simulation studies and real data analyses, we demonstrated that pruning provides an improvement over simple 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.

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

1 **Introduction**

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

17

18 A well-known polygenic approach for genotype data is polygenic risk scores (PRS),
19 which are weighted sum of risk alleles of a pre-selected number of genetic variants¹¹.
20 Recently, many efforts have focused on transferring PRS approaches to DNA
21 methylation data to construct methylation risk scores (MRS), which are defined as

22 weighted sums of the individuals' DNAm values of a pre-selected number of CpGs¹⁰.
23 However, there are many methodological challenges in constructing DNA
24 methylation risk scores^{10,12,13}. One of the problems is that DNAm is influenced by
25 ancestry, which captures genetic ancestry (differences in the genome related to
26 ancestry) as well as social determinants of health such as racism and discrimination,
27 socioeconomic status, and environmental effects¹⁴. Thus, ideally, when external
28 weights are used for the calculation of MRS, these weights should be assessed in a
29 population with the same ancestry as the study samples. However, current
30 epigenetic literature remains limited by the lack of diversity, with most focusing on
31 European populations¹⁵, therefore making it difficult to identify appropriate weights
32 for MRS for other populations. While it is well known that PRSs are not applicable
33 across different ancestries^{16,17}, little is known about the performance of MRS across
34 multi-ancestry populations.
35
36 Currently, there are two popular approaches to construct MRS. The first one is to
37 use penalized regression models such as Elastic Net and LASSO
38 regularization^{10,18,19}, which usually requires individual-level DNAm data. When only
39 summary-level statistics are available, individual CpGs that reached genome-wide
40 significance in an external epigenome-wide association studies (EWAS) are selected
41 and the beta-coefficients estimated from EWAS are used as weights to generate
42 MRS¹⁰. However, research in PRS has shown that the optimal p-value threshold

43 strongly depends on the data²⁰, and including a larger proportion of variants could
44 potentially capture more of the phenotype variation²¹. Moreover, most MRS from the
45 second approach do not consider DNA co-methylation, defined as proximal CpGs
46 with correlated DNAm across individuals²², which could potentially bias the
47 generation of MRS. Shah et al 2015 proposed to remove redundant CpGs by
48 keeping the most significant CpGs in co-methylation^{23,24}, however, a window to
49 define DNA co-methylation needs to be pre-defined and the effect of accounting for
50 DNA co-methylation was not evaluated.

51

52 One of the most widely used PRS approaches to deal with single nucleotide
53 polymorphisms (SNPs) in high linkage disequilibrium (LD) and to identify p-value
54 thresholds with the best prediction accuracy is the pruning and threshold (P+T)
55 method²⁵. In the P+T approach, the correlation square (R^2) for SNPs within a close
56 genetic distance is calculated and less significant SNPs that are correlated with an
57 R^2 greater than a particular value (LD pruning)²⁶ are removed. Next, several p-value
58 thresholds are tested to maximize the prediction accuracy of the derived PRS (p-
59 value thresholding)^{26,27}. Theoretically, the P + T approach could be applied to
60 generate MRS, however, there is no standard procedure on how to conduct pruning
61 for DNAm data and the performance of such MRS across multi-ancestry populations
62 remains unknown.

63

64 Here, we propose to use the Co-Methylation with genomic CpG Background
65 (CoMeBack), a tool that uses a sliding window to estimate DNA co-methylation, to
66 account for correlations of DNAm at proximal CpG sites²², and pair it with p-value
67 thresholding to construct P+T CoMeBack MRS. CoMeBack uses unmeasured
68 intermittent CpGs from the human reference genome to link array probes in hope of
69 reducing false positives while improving the identification of biologically relevant co-
70 methylation²². We conducted simulation studies based on data from an adult
71 population consisting of three groups of different ancestries (Indian, White and Black,
72 n = 1,199) to evaluate the prediction performance of P+T CoMeBack MRS and how
73 it changes across multi-ancestry population. Next, we applied the P+T CoMeBack
74 approach to DNAm data from the Drakenstein Child Health Study (n=270)²⁸, a multi-
75 ancestry birth cohort from South Africa, to evaluate the performance of MRS for
76 maternal smoking status. Our simulation study and real data application
77 demonstrated that the P+T approach improves the predictive accuracy of MRS over
78 methods that do not account for co-methylation and has similar performance as
79 LASSO regression, which requires access to the raw DNAm data. We also showed
80 that MRS built upon the data from a population of one genetic ancestry could
81 achieve high prediction performance among populations of other genetic ancestries,
82 but the performance might differ in the presence of environmental/cultural/social
83 differences associated with ancestry.

84

85 **Materials and methods**

86 **P+T CoMeBack approach for MRS**

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

94
95 In our P+T CoMeBack method, co-methylation pruning is completed by applying
96 CoMeBack to DNAm data of the testing dataset or a reference panel²². Specifically,
97 CoMeBack chains two adjacent array probes if the following requirements are met:
98 1. two probes are less than 2kb apart; 2. the reference human genome annotation
99 shows a set of unmeasured genomic CpGs between them; 3. the density of
100 unmeasured genomic CpGs between them is at least one CpG every 400bp.

101 Chaining of adjacent array probes continues until an array probe does not meet the
102 requirements, which will form a unit where multiple CpGs are chained together.

103 Correlations between DNAm levels will then be calculated for all array probes inside
104 each unit. If all pairs of adjacent probes in a unit have a correlation square (R^2)
105 greater than 0.3, such unit will be declared as a co-methylated region (CMR).

106 Pruning is conducted by only keeping one CpG site per CMR in the dataset, the one
107 with the lowest (most significant) P-value in the EWAS summary statistics.

108

109 P + T CoMeBack will be compared to the standard pruning approach, in which less
110 significant SNPs that are correlated with an $R^2 > 0.3$ and located within 2000bp of
111 each other are being removed.

112

113 Next, P-value thresholding step (T) is performed for the pruned set of CpG sites.

114 Specifically, the P-value thresholding step (T) is performed by applying different P-

115 value thresholds (e.g., P-value thresholds $\in [0.05, 0.005, 5 \times 10^{-4}, 5 \times 10^{-5}, \dots]$) and

116 only including those CpG sites in the final MRS calculation that reached a P-value

117 below those thresholds in the EWAS summary statistics.

118

119 Finally, for each P-value threshold, MRS are calculated as a weighted sum of DNAm

120 β values (β value = methylated allele intensity / (unmethylated allele intensity +

121 methylated allele intensity + 100), ranging from 0 representing unmethylation to 1 for

122 complete methylation) of the selected CpGs, where the weights are the

123 corresponding effect sizes for each CpG from the EWAS summary statistics. The

124 squared correlation (R^2) between the phenotype of interest and MRS obtained using

125 each P-value threshold is calculated to represent the prediction accuracy. The P-

126 value threshold that produces MRS with the highest prediction accuracy in the

127 testing data set is selected as the optimal P-value and the corresponding MRS is
128 used for downstream analysis. The pipeline for generating P+T MRS is written in an
129 R script, which is available at GitHub ([https://github.com/jche453/Pruning-
130 Thresholding-MRS.git](https://github.com/jche453/Pruning-Thresholding-MRS.git)).

131
132 In our simulation studies and real data application, we compare the P+T CoMeBack
133 MRS approach to the standard P+T and T approach, which refers to an approach in
134 which the MRS is calculated by only thresholding, not accounting for correlations
135 between included CpGs (no pruning).

136

137 **Simulation studies**

138 To validate the performance of the proposed P+T MRS approach, we conducted
139 simulation studies based on whole blood Illumina Infinium Human Methylation 450K
140 BeadChip data from an ethnically heterogeneous discovery cohort composed of
141 several publicly available datasets (GSE55763, GSE84727, GSE80417, GSE111629
142 and GSE72680)²². Intra-dataset normalization and batch effects correction were
143 performed using ComeBat function in R-package sva²⁹, followed by merging of
144 datasets and correction for inter-dataset batch effects using the same function. After
145 the removing XY chromosome binding, non-CpG, cross-hybridizing probes and
146 probes that are in close distance with common SNPs, there were 386,362 CpGs left

147 for MRS analysis. We randomly selected 1,199 adults (898 Indians, 136 Blacks and
148 165 Whites) to conduct the simulation studies.

149
150 CoMeBack was applied to the DNA methylation β values of the 386,362 CpGs to
151 obtain CMR. In each simulation, 10 of the 386,362 CpGs were randomly selected to
152 be causal, $k\%$ ($k = 30, 50, 70$ or 100) of which are in a CMR with other CpGs. At
153 most, one CpG would be causal in each CMR.

154
155 The causal CpGs were randomly assigned a “true” effect size from a uniform
156 distribution as $w_i \sim U(-0.5, 0.5)$. We then simulated a phenotype for the j -th subject
157 as follow:

$$158 \quad Y_j = \sum_{i=1}^{10} w_i m_{ij} + \varepsilon_j, \varepsilon_j \sim N(0, \delta^2),$$

159 where m_{ij} is the DNAm β value of causal CpG site i of the j -th subject, and ε_j is an
160 error term that follows a normal distribution. Different δ^2 were set to ensure that the
161 targeted variance of phenotype explained by DNAm alone equals 10%, 30% 50% or
162 80%.

163
164 We also simulated a second phenotype Y_j^* for the j -th subject, which was directly
165 affected by ancestry using European ancestry as reference:

$$166 \quad Y_j^* = \sum_{i=1}^{10} w_i m_{ij} + a * (\text{if Indian}) + b * (\text{if Black}) + \varepsilon_j^*, \varepsilon_j^* \sim N(0, \delta^{*2})$$

167 Effect of ancestry in our simulations is simulated as the effect of genetic ancestry
168 assuming there were no complex social determinants involved in the causal
169 pathway. Different δ^{*2} were used so that the variance of phenotype that was
170 explained by DNAm and ancestry together equals 20%, 50% or 80. For our
171 simulations, effect a was set to 0.1 and b to 0.2. In each simulation, both simulated
172 phenotypes Y_j and Y_j^* share the same epigenetic liability ($\sum_{i=1}^{10} w_i m_{ij}$).
173
174 In each simulation, for fair comparison, 762 Indians were randomly chosen as the
175 training dataset so that there were at least 136 people left for each race group in the
176 testing dataset. Associations between CpGs and each of the two simulated
177 phenotypes were assessed by robust linear regression model using limma R
178 package³⁰ in the training dataset. We calculated top 10 principal components (PCs)
179 from DNAm of 386,362 CpGs³¹ and used EpiDISH to estimate cell type proportions
180 of each CpGs³². We observed that in our simulation dataset, top 10 PCs are highly
181 correlated with cell type proportions (Supplement figure 1A), and using either
182 summary statistics adjustment for top 10 PCs or summary statistics adjusted for cell
183 type proportions would lead to almost identical prediction performance of MRS
184 (Supplement figure 1B). Thus, to account for population stratification and cell type
185 difference, we adjusted for the top 10 PCs in our main analyses. The summary
186 statistics (effect size and P-values) obtained from association tests in the training
187 data were saved and later used to construct MRS in the testing dataset. We

188 repeated 1000 simulations per scenario to evaluate the prediction accuracy (R^2),
189 power and type 1 error rate of the P+T MRS. Linear regression analysis was used to
190 access the association between MRS and simulated phenotypes, and power is
191 defined as the proportion of simulations where MRS were significantly associated
192 with the simulated phenotype with at α level of 0.05. To estimate type 1 error rate,
193 we first obtained a null association between MRS values and simulated phenotype
194 values by permutation of MRS values. Linear regression analysis was used to
195 access the association between permuted MRS and simulated phenotypes, and
196 type 1 error rate is defined as the proportion of simulations where permuted MRS
197 were significantly associated with the simulated phenotype.

198

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

207

208 **Application study of smoking MRS**

209 To evaluate the performance of the P+T CoMeBack approach in a real data setting,
210 we applied the P+T CoMeBack approach to calculate a MRS for maternal smoking
211 status during pregnancy using cord blood DNAm data from newborns in the South
212 African Drakenstein Child Health Study (DCHS), a multi-ancestry longitudinal study
213 investigating determinants of early child development³³. There were 145 Black
214 African infants and 115 Mixed ancestry infants in the DCHS. A detailed description of
215 the enrollment process, inclusion criteria, variables measurement and ethical
216 approval of the study have been previously published^{33,34}.

217

218 Cotinine levels were measured in urine provided by mothers within four weeks of
219 enrollment and classified as <499 ng/ml (non-smoker), or ≥500 ng/ml (active
220 smoker)²⁸. Cord blood was collected at time of delivery and used to measure DNA
221 methylation by either MethylationEPIC BeadChips (EPIC, n=145) or the Illumina
222 Infinium HumanMethylation450 BeadChips (450K, n=103)^{33,34}, followed by quality
223 control and normalization to calculate β values (details have been published
224 elsewhere)³⁵.

225

226 Summary statistics for the calculation of MRS were obtained from a study that meta-
227 analyzed the associations between newborn blood DNA methylation and sustained
228 maternal smoking during pregnancy among 5,648 mother-child pairs as part of the
229 Pregnancy and Childhood Epigenetics (PACE) Consortium (Table 1)³⁶. The

230 participants of all cohorts used in the meta-analysis except one were of European
231 ancestry.

232

233 In addition, we compared P+T CoMeBack MRS to three previously published MRS
234 for maternal smoking during pregnancy (Reese MRS, Richmond 19 MRS, Richmond
235 568 MRS; Table 1). Reese MRS model was trained among 1,068 newborns of
236 European ancestry in the Norwegian Mother and Child Cohort Study, while
237 Richmond 568 MRS and Richmond 19 MRS was trained in multi-ancestry newborns
238 (N=6,685) and children around 6.8 years old (N=3,187) in PACE Consortium
239 respectively. The training population for Reese MRS and Richmond 19 MRS
240 overlapped with the training population for summary statistics used in P+T
241 CoMeBack MRS in our study. Reese et al. used a LASSO regression to select CpGs
242 for Reese MRS, which is a weighted sum of DNAm β values of 28 CpGs with
243 weights estimated from the LASSO regression³⁷. Richmond 19 MRS is a weighted
244 sum of DNAm β values of 19 CpGs that were significantly associated with prenatal
245 smoking in an EWAS conducted in peripheral blood from children of averaged 6.8
246 years age (Richmond 19 MRS)^{38,39}. In the same study, Richmond 568 MRS was
247 proposed based on 568 CpGs that were significantly associated with prenatal
248 smoking in cord blood^{38,39}. We obtained the weights of reported CpGs from the
249 mentioned studies and applied them to DNAm data in DCHS to generate Reese
250 MRS, Richmond 19 MRS and Richmond 568 MRS.

251

252 Linear regressions were used to assess the associations between maternal smoking
253 status and each MRS, controlling for ancestry (in pooled samples), cell type
254 proportions and top 5 PCs calculated from genotypes. In order to obtain comparable
255 beta-coefficients and standard errors across different MRS, each MRS was divided
256 by their interquartile range (IQR) before linear regression analysis.

257

258 **Results**

259 **Simulation results**

260 We compared the prediction performance of P+T CoMeBack MRS to the T method
261 among 136 Indians in the test data across different simulation scenarios (Figure 1).
262 Figure 1A shows that P+T CoMeBack MRS that account for co-methylation between
263 CpGs have stable prediction performance when proportion of causal CpGs located in
264 a CMR (k%) varies. P+T without CoMeBack had similar prediction performance while
265 the T method had a slightly lower prediction performance. While the P+T CoMeBack
266 MRS showed subtle improvement over T method when V_{DNAm}^2 is 80% (Figure 1A),
267 the difference between P+T CoMeBack MRS and the T method decreases as the
268 V_{DNAm}^2 decreases. This is likely because as variance explained by DNAm decreases,
269 there is less power for association testing, and it becomes increasingly difficult to
270 distinguish real signals from statistical noise while generating the summary statistics.

271

272 Next, we assessed the performance of P+T CoMeBack, P+T and T method across
273 different ancestries and among multi-ancestry populations. All three methods
274 achieved a high power (> 95%) and a low type 1 error rate (~ 5%) within each
275 ancestry in most scenarios for both phenotypes except when the phenotype variance
276 explained by DNA methylation is 10% or 30% (Supplement table 1-4).

277

278 Whether the simulated phenotypes were independent of ancestry or not, MRS
279 among Whites and Blacks achieved a prediction R^2 as high as among Indians, which
280 should have the best prediction of the simulated phenotypes since weights were
281 obtained from Indian training samples (Figure 2). Findings were similar when the
282 phenotype variance explained by DNA methylation was reduced from 80% to 10%,
283 30% or 50% (Supplement Figure 2). When the phenotypes are not associated with
284 ancestry (Figure 2A), the three MRS analyses strategies (stratification, joint analysis
285 and standardization) lead to nearly identical results. However, when the phenotypes
286 are ancestry-dependent, both joint analysis and standardization of MRS showed very
287 poor prediction of the phenotypes (Figure 2B).

288

289 **MRS of maternal smoking status**

290 Figure 3 shows the prediction performance of MRS for maternal smoking status
291 among DCHS newborns. As the p-value threshold decreases, the prediction
292 accuracy of the resulting MRS increases before reaching a plateau, demonstrating

293 the importance of P-value thresholding in MRS to control for noise. Among mixed
294 ancestry newborns, P+T CoMeBack MRS of smoking status excluded 22 CpGs in
295 pruning and achieved a prediction R^2 of 29.5% using P-value threshold of 5×10^{-22} ,
296 while the standard P+T without CoMeBack had a lower prediction R^2 of 26.2% using
297 P-value threshold of 5×10^{-10} and the best T method MRS had the lowest prediction
298 accuracy (24.5%) using P-value threshold 5×10^{-9} , confirming the benefits of pruning
299 in MRS calculation (Figure 3A). All three MRS had lower prediction performance for
300 maternal smoking among Black African infants (10.9%, and 8.0% respectively)
301 (Figure 3B), which is likely due to the low prevalence of smokers among mothers of
302 Black African infants in DCHS (13%) compared to mothers of mixed ancestry infants
303 (49%) (Supplement Figure 3). Additionally, the distributions of all MRS in Black
304 African infants and mixed ancestry infants were similar within each category of
305 maternal smoking status (Supplement Figure 4), confirming that the difference of
306 prediction R^2 between Mixed and Black infants is less likely due to ancestry-related
307 factors other than prevalence of maternal smoking. Joint-analysis of P+T CoMeBack
308 MRS showed a prediction accuracy of 20.4%, which is between the prediction
309 accuracy of P+T CoMeBack MRS among Black African infants and mixed ancestry
310 infants (Figure 3C). Standardization approach did not improve the performance of
311 MRS (Figure 3D).
312

313 We next compared the prediction accuracy and distribution of P+T CoMeBack MRS
314 to other established MRS for maternal smoking during pregnancy and newborn
315 DNAm (Figure 4). Overall, P+T CoMeBack and Reese MRS had stable and similar
316 classification performance in all analyses compared to other MRS. P+T CoMeBack
317 MRS and Reese MRS showed a similar prediction R^2 among both Black and Mixed
318 ancestry infants, which are better than other smoking MRS (Figure 4A). P+T
319 CoMeBack MRS had the largest AUC (0.820) in the ROC curve among mixed infants
320 (Figure 4B) but a smaller AUC than Reese MRS in Black infants and joint-analysis
321 (Figure 4C-D). Further, all 6 MRS showed significant association with smoking status
322 in Mixed infants, Black infants and joint-analysis (Table 2), showing the promise of
323 using MRS to capture the overall DNAm signals in association testing.

324

325 **Discussion**

326 Based on the well-established P+T CoMeBack framework in PRS, we developed
327 P+T CoMeBack MRS, which aggregates EWAS signals and could potentially be
328 used as a biomarker in association studies where single CpGs do not achieve
329 significance^{4,40,41}. The proposed P+T CoMeBack MRS approach uses CoMeBack for
330 co-methylation pruning and evaluates multiple P-value thresholds to maximize
331 prediction performance. Such MRS could potentially serve as a powerful dimension
332 reduction approach for mediation and multi-omics integration analyses^{4,40-43} as well
333 as biomarkers of individual disease risk in a clinical setting⁴⁴⁻⁴⁶.

334

335 Overall, our simulation studies demonstrated good performance of P+T CoMeBack
336 MRS for predicting phenotypes of interest with good statistical power and well-
337 controlled type 1 error. We demonstrated that the prediction accuracy of MRS
338 reflects the variance of phenotype that is explained by DNAm. By accounting for
339 inter-correlation between CpGs, P+T CoMeBack MRS and P+T without CoMeBack
340 showed a slightly better performance than the standard T method. In the real data
341 application, we observed the best prediction of maternal smoking status when using
342 P+T CoMeBack, which confirms the usefulness of accounting for co-methylation and
343 demonstrates the ability of CoMeBack to control for false discover of CMR and
344 usefulness in constructing MRS²². However, we note that P+T CoMeBack MRS
345 could still have poor prediction performance if the external EWAS is underpowered
346 or subject to bias.

347

348 In the prediction of maternal smoking status, P+T CoMeBack MRS showed
349 comparable performance to Reese MRS, which was derived using the LASSO
350 method³⁷. When predictors are highly correlated, LASSO typically selects one of the
351 correlated predictors and shrinks the effect size of the rest to zero, which might
352 produce similar results to our pruning procedure in developing MRS. One of the
353 advantages of P+T CoMeBack MRS is that it is based on EWAS summary statistics
354 which are often publicly available, hence making it a valuable approach, as it is often

355 difficult to obtain individual DNAm data from an external cohort. Additionally, P+T
356 CoMeBack MRS can make use of meta-analysis-type summary statistics, which
357 aggregates results from multiple studies to improve association estimates. In
358 contrast, to construct MRS like Reese MRS, individual DNAm data are usually
359 required to perform a LASSO regression, and these are often not accessible.
360 Recently, novel penalized regressions have been proposed to generate PRS with
361 only GWAS summary statistics and publicly available reference data⁴⁷, but their
362 applications to EWAS summary statistics for MRS have not been investigated. To
363 develop MRS for different exposures and outcomes, we urge EWAS studies to make
364 their genome-wide summary statistics publicly available.

365
366 In our simulation studies, weights obtained from Indian training samples were
367 applied to generate P+T CoMeBack MRS, thus MRS among Indian testing samples
368 were assumed to have the best prediction of the simulated phenotypes. However,
369 MRS among Whites and Blacks also achieved a prediction accuracy as high as
370 among Indians for both simulated phenotypes suggesting that genetic ancestry does
371 not contribute to difference in prediction abilities of MRS across multi-ancestry
372 population. This is likely because we assumed all ancestries share the same causal
373 CpGs and effect sizes. However, in the real world, this assumption could possibly be
374 violated for many phenotypes. Unlike ancestry in our simulation studies, ancestry in
375 the real world is complex. The meaning of ancestry could be different in different

376 regions/nations, and “effect of ancestry” involves the joint effects of ancestry-
377 associated social determinants of health and environmental effects , and cultural
378 context⁴⁸. Ancestry, along with environment and social differences associated with it,
379 could affect both MRS and phenotypes in numerous causal pathways and potentially
380 modify the effect of MRS on the phenotypes. Thus, even if all ancestries indeed
381 share the same causal CpGs and effect sizes, it might still not be sufficient to
382 disentangle the relationship between ancestry, DNAm and phenotype of interest.
383 This may greatly impact the transferability of MRS across different ancestries, which
384 could be the reason why we observed an inconsistency of performance of P+T
385 CoMeBack MRS in terms of their distributions and predictions across multi-ancestry
386 population in the real data analyses. In practice, we recommend that researchers
387 conduct MRS analyses stratified by ancestry first and evaluate the effect of ancestry
388 on MRS analyses before pooling participants together for a joint analysis.
389
390 In our real data application, summary statistics for smoking were obtained from a
391 cohort with mainly people of European ancestry³⁶. MRS of smoking among mixed
392 ancestry infants achieved a prediction accuracy of nearly 30%. However, the
393 prediction accuracy of P+T CoMeBack MRS among Black African infants was only
394 10.9%. We suspect that the difference was largely due to the prevalence of active
395 smoking among mothers of Black African infants being lower than those of mixed

396 ancestry infants (13% vs 49%), which is similar to how the prevalence of outcome
397 affects the predictive ability of PRS⁴⁹.

398

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

413

414 In conclusion, P+T in general and P+T using CoMeBack in particular, provides an
415 improvement for prediction of phenotype of interest, over T method that does not
416 account for co-methylation between CpGs. In contrast to PRS, using existing

417 summary statistics that were derived from European populations can be used to
418 calculate MRS in other ancestries, thus reducing the ancestry/ethnicity disparity in
419 medical research. However, caution is needed in the analyses and interpretation of
420 MRS results across multi-ancestry populations. More investigations of MRS are
421 urged to further improve their prediction accuracy and translational values, also in
422 combination with other clinical and non-clinical variables, especially among multi-
423 ancestry population. With the current increase of large consortia-led EWAS for
424 different exposures and health outcomes (e.g., the PACE consortium), we believe
425 the predictive performance of MRS will continue to increase, and the P+T CoMeBack
426 method has the potential to be widely used for risk prediction and association testing.

427

428 **Acknowledgments:**

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

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

MRS	Training dataset publication	Training Population	Phenotype	MRS publication	P-value threshold/ Method	No. of CpG sites (joint-analysis)
P+T CoMeBack MRS	Sikdar et al. 2019 ³⁶	Multi-ethnic newborns (mainly White, N=5,648)	Most cohorts ascertained sustained smoking during pregnancy by questionnaires; two cohorts incorporated cotinine-based smoking measure	-	5x10 ⁻²² (Mixed)	21 (43 passed P-value threshold and 22 excluded by pruning)
					5x10 ⁻²⁴ (Black)	20 (42 passed P-value threshold and 22 excluded by pruning)
					5x10 ⁻²² (Pooled)	21 (43 passed P-value threshold and 22 excluded by pruning)
P+T MRS	Sikdar et al. 2019 ³⁶	Multi-ethnic newborns (mainly White, N=5,648)	Same as above	-	5x10 ⁻¹⁰ (Mixed)	198 (233 passed P-value threshold and 35 excluded by pruning)
					5x10 ⁻³⁶ (Black)	4 (26 passed P-value threshold and 22 excluded by pruning)
					5x10 ⁻²⁴ (Pooled)	8 (42 passed P-value threshold and 34 excluded by pruning)

					5x10 ⁻⁹ (Mixed)	344
T MRS	Sikdar et al. 2019 ³⁶	Multi-ethnic newborns (mainly White, N=5,648)	Same as above	-	5x10 ⁻²⁴ (Black)	42
					5x10 ⁻¹⁶ (Pooled)	72
Reese MRS	Reese et al. 2017 ³⁷	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 ³⁷	Logistic LASSO regression	28
Richmond 568 MRS	Joubert et al. 2016 ³⁹	Multi-ethnic newborns (N=6,685)*	Maternal smoking during pregnancy via questionnaires	Richmond et al. 2018 ³⁸	Robust linear regression; Bonferroni corrected P-value < 0.05	568
Richmond 19 MRS	Joubert et al. 2016 ³⁹	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.

Table 2. Association between maternal smoking status and MRS in DCHS.

MRS	Mixed			Black			Pooled (Joint-analysis)		
	Beta-coefficient*	Standard Error	P-value	Beta-coefficient	Standard Error	P-value	Beta-coefficient	Standard Error	P-value
P+T CoMeBack MRS	0.88	0.12	5.65 x10 ⁻¹¹	0.76	0.18	5.74 x10 ⁻⁵	0.82	0.10	3.01x10 ⁻¹⁴
P+T MRS	0.71	0.11	1.05 x10 ⁻⁸	0.52	0.14	3.29 x10 ⁻⁴	0.69	0.10	1.32x10 ⁻¹⁰
T MRS	0.80	0.13	7.35 x10 ⁻⁹	0.70	0.19	2.82 x10 ⁻⁴	0.64	0.09	2.55 x10 ⁻¹¹
Reese MRS	0.97	0.13	3.77 x10 ⁻¹¹	0.73	0.16	2.06 x10 ⁻⁵	0.77	0.09	2.50 x10 ⁻¹⁵
Richmond 568 MRS	0.85	0.14	9.84 x10 ⁻⁹	0.54	0.16	7.82 x10 ⁻⁴	0.65	0.10	1.57 x10 ⁻¹⁰
Richmond 19 MRS	0.76	0.14	5.37E x10 ⁻⁷	0.70	0.19	3.78 x10 ⁻⁴	0.70	0.11	2.05 x10 ⁻⁹

***Beta-coefficients indicate change in interquartile range (IQR) of MRS between smokers and non-smokers.**

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. Horvath S. DNA methylation age of human tissues and cell types. *Genome biology*. 10/21 06/10/received 10/21/accepted 2013;14(10):R115-R115. doi:10.1186/gb-2013-14-10-r115
19. Thompson M, Hill BL, Rakocz N, et al. Methylation risk scores are associated with a collection of phenotypes within electronic health record systems. *npj Genomic Medicine*. 2022/08/25 2022;7(1):50. doi:10.1038/s41525-022-00320-1
20. 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
21. 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
22. 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
23. Shah S, Bonder MJ, Marioni RE, et al. Improving Phenotypic Prediction by Combining Genetic and Epigenetic Associations. *Am J Hum Genet*. Jul 2 2015;97(1):75-85. doi:10.1016/j.ajhg.2015.05.014
24. Odintsova VV, Rebattu V, Hagenbeek FA, et al. Predicting Complex Traits and Exposures From Polygenic Scores and Blood and Buccal DNA Methylation Profiles. *Front Psychiatry*. 2021;12:688464. doi:10.3389/fpsyt.2021.688464
25. 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
26. 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>
27. Choi SW, O'Reilly PF. PRSice-2: Polygenic Risk Score software for biobank-scale data. *GigaScience*. 2019;8(7)doi:10.1093/gigascience/giz082
28. 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
29. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*. Mar 15 2012;28(6):882-3. doi:10.1093/bioinformatics/bts034

30. 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
31. Barfield RT, Almlı 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
32. 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
33. 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
34. 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
35. 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
36. 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
37. 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
38. 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
39. 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>
40. 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
41. 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
42. 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
43. 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

44. 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
45. 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
46. 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
47. 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
48. 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
49. 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

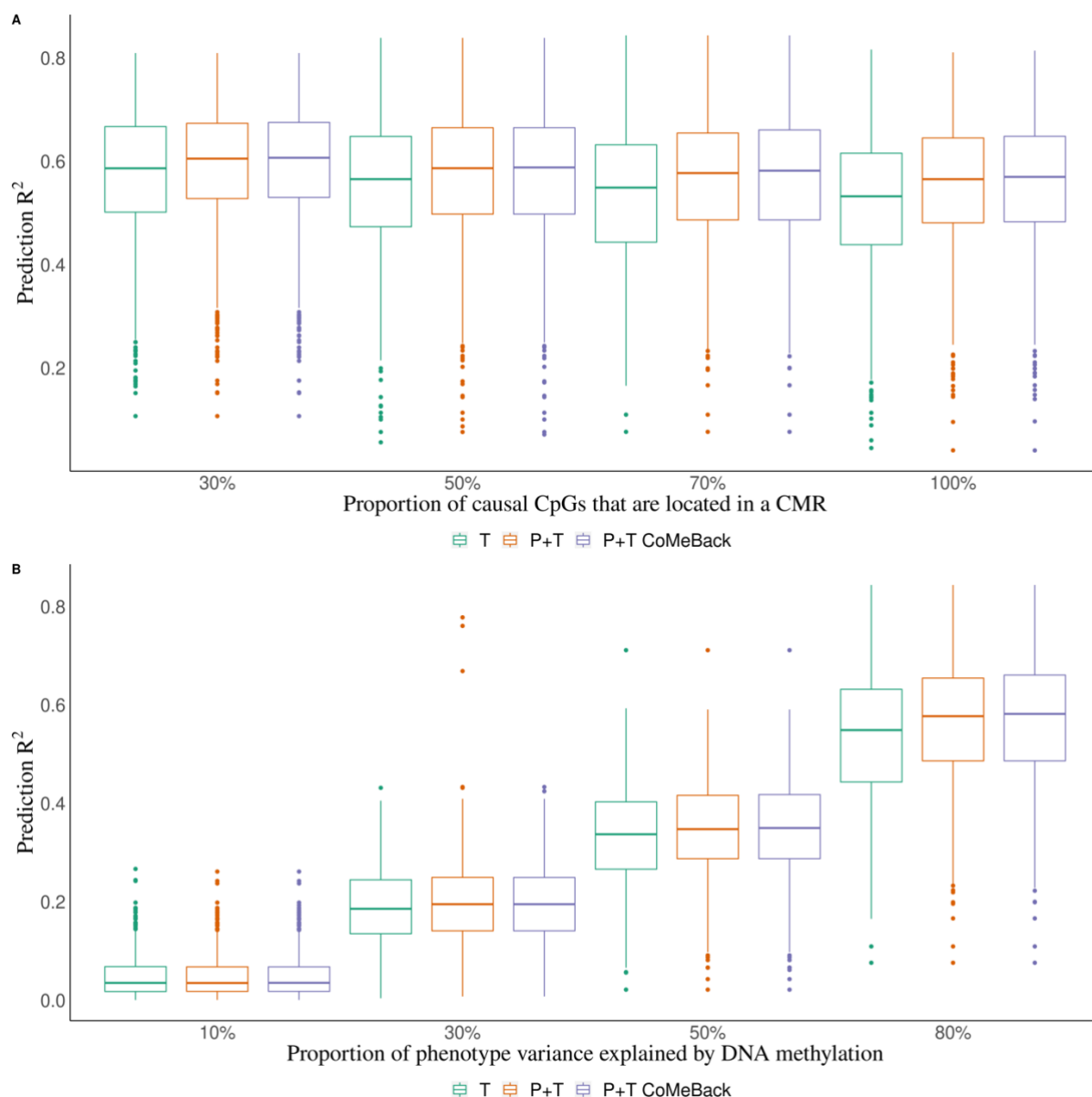


Figure 1. Simulation study. Prediction R^2 of P+T CoMeBack, P+T and T method in dependence of (A) the proportion of causal CpG sites in CMRs and (B) proportion of phenotype variance explained by DNA methylation, among Indian participants. For each simulation, the discovery cohort was repeatedly and randomly 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.

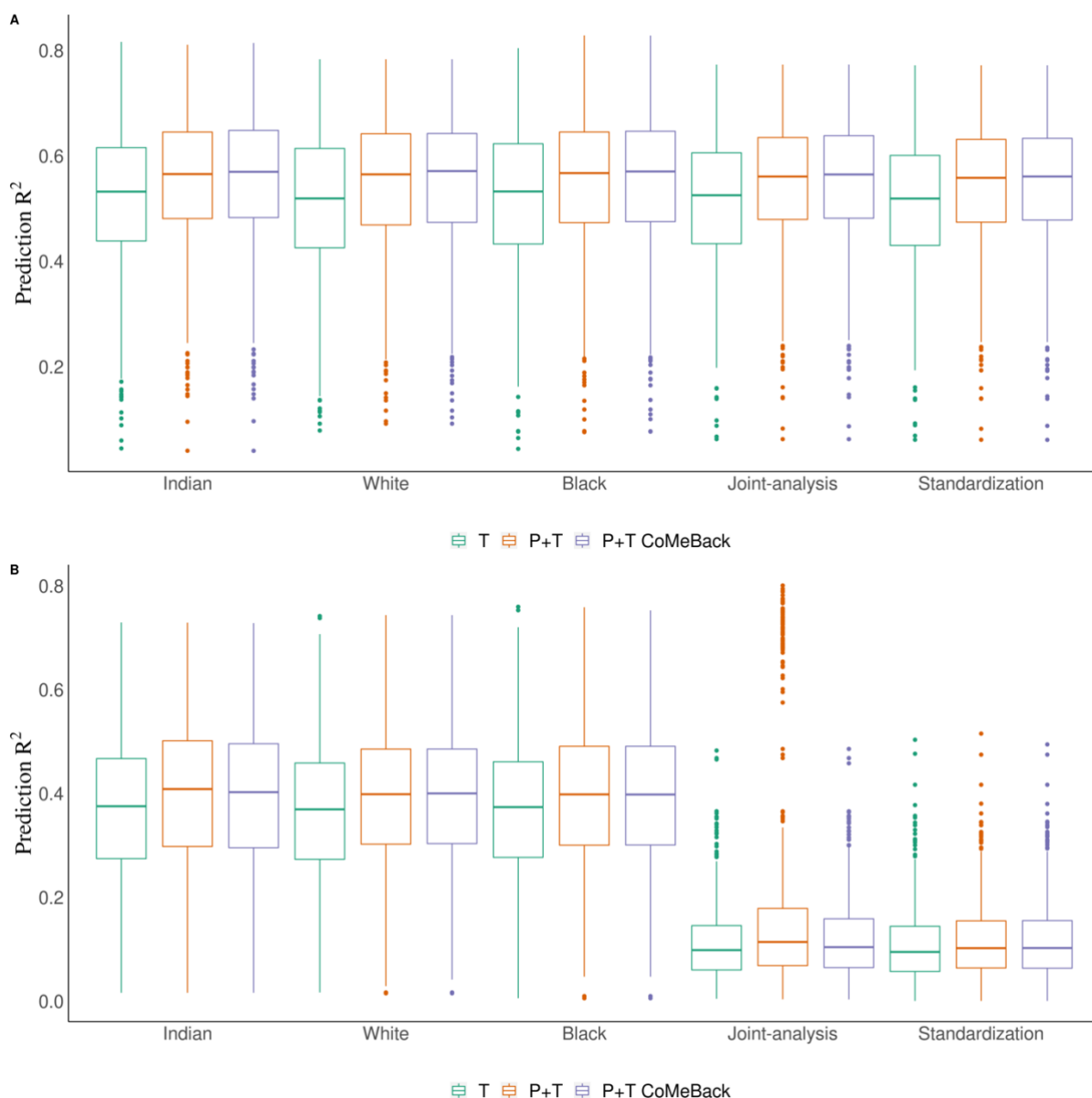


Figure 2. Simulation study. Prediction R^2 of P+T CoMeBack and T approach across different racial groups and among multi-ancestry populations. For each simulation, the discovery cohort was repeatedly and randomly split into a training set comprising 762 Indians and a testing set comprising 136 people 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.

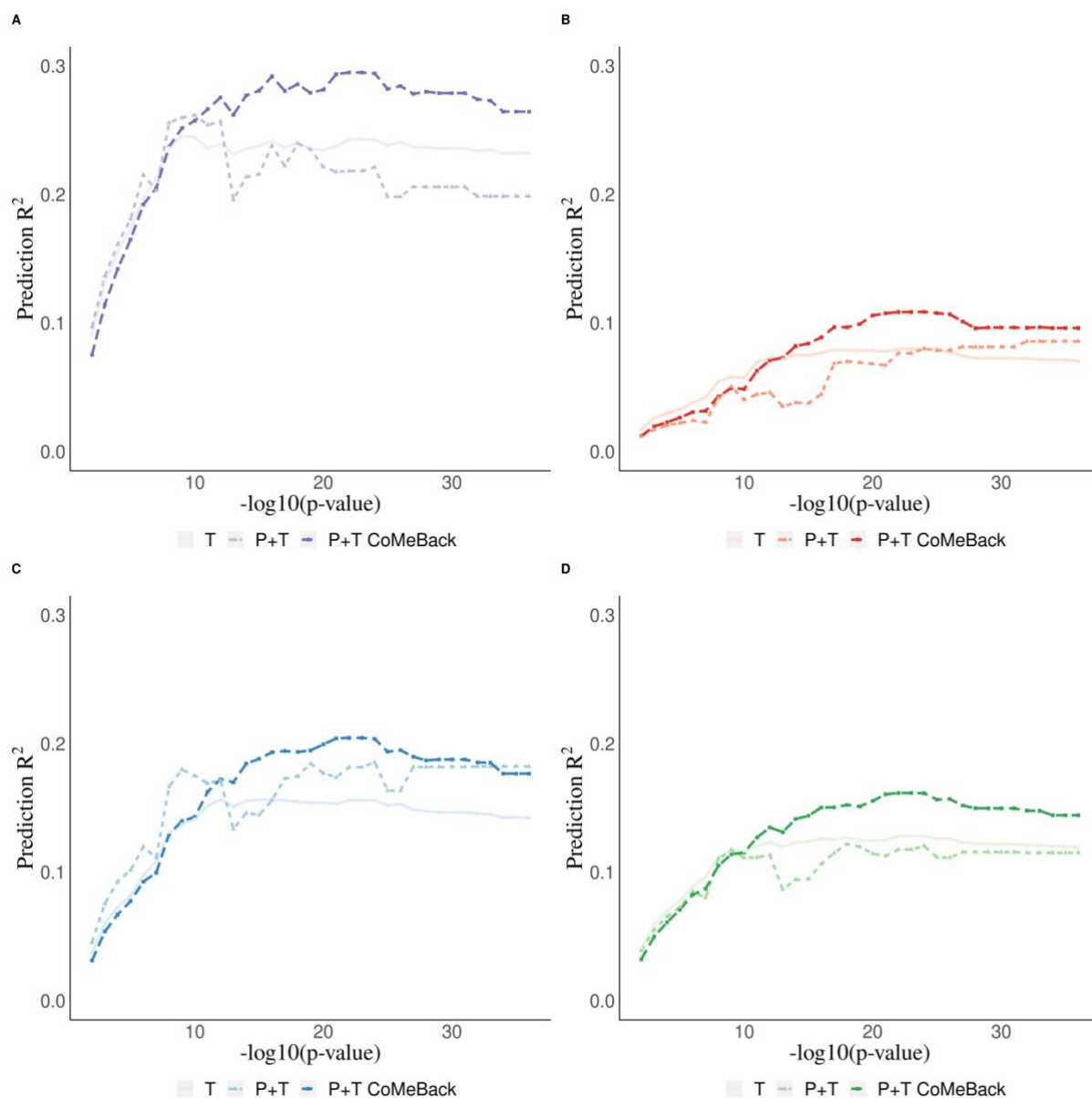


Figure 3. Real data application. MRS for the prediction of maternal smoking during pregnancy using cord blood DNA methylation data from newborns in the South African Drakenstein Child Health Study (DCHS). Prediction R^2 of maternal smoking status is shown stratified for **A.** Mixed infants. **B.** Black infants. **C.** joint-analysis (all subjects pooled from all ancestries) **D.** Standardization (standardizing MRS within each ancestry and merging all subjects before analyses)

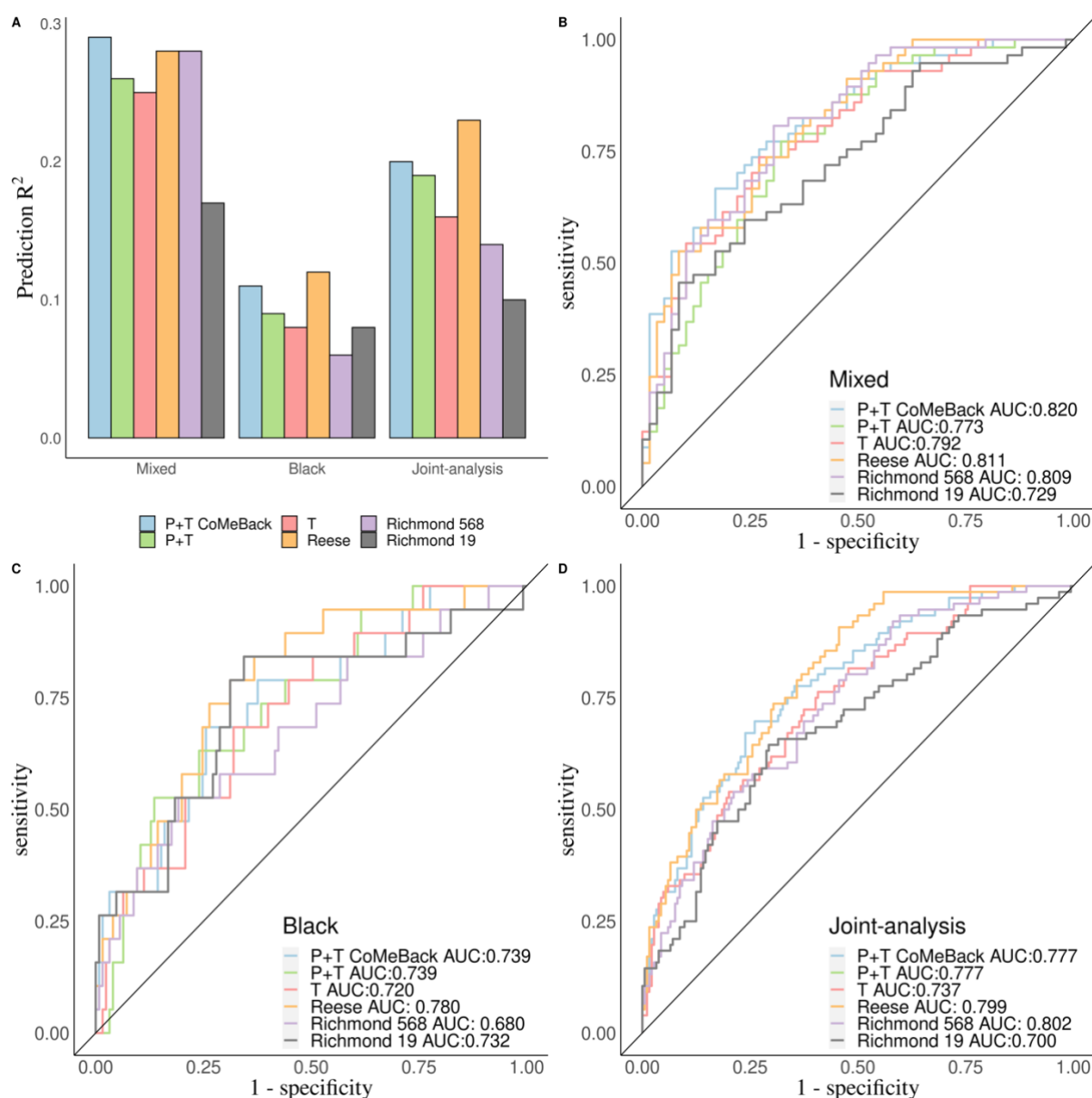


Figure 4. Real data application. Comparison of P+T CoMeBack method to P+T, T and 3 other published MRS for predicting maternal smoking status in the South African Drakenstein Child Health Study (DCHS). A Prediction R² of all 6 MRS methods for Mixed infants, Black infants and pooled samples (joint-analysis). A receiver operating characteristic (ROC) curve comparing prediction performance of all 6 MRS among **(B)** Mixed infants, **(C)** Black infants and **(D)** pooled samples (joint-analysis).