

# 1           **Methods for Mediation Analysis with High-Dimensional DNA** 2           **Methylation Data: Possible Choices and Comparison**

3  
4           Dylan Clark-Boucher<sup>1</sup>, Xiang Zhou<sup>2</sup>, Jiacong Du<sup>2</sup>, Yongmei Liu<sup>3</sup>, Belinda L  
5           Needham<sup>4</sup>, Jennifer A Smith<sup>4,5</sup>, Bhramar Mukherjee<sup>2,4</sup>

6  
7           <sup>1</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA

8           <sup>2</sup>Department of Biostatistics, University of Michigan, Ann Arbor, MI

9           <sup>3</sup>Department of Medicine, Divisions of Cardiology and Neurology, Duke University Medical  
10           Center, Durham, NC

11           <sup>4</sup>Department of Epidemiology, University of Michigan, Ann Arbor, MI

12           <sup>5</sup>Survey Research Center, Institute for Social Research, University of Michigan, Ann Arbor, MI

## 13 14           **Abstract**

15           Epigenetic researchers often evaluate DNA methylation as a mediator between social/environmental  
16           exposures and disease, but modern statistical methods for jointly evaluating many mediators have not  
17           been widely adopted. We compare seven methods for high-dimensional mediation analysis with  
18           continuous outcomes through both diverse simulations and analysis of DNAm data from a large national  
19           cohort in the United States, while providing an R package for their implementation. Among the  
20           considered choices, the best-performing methods for detecting active mediators in simulations are the  
21           Bayesian sparse linear mixed model by Song et al. (2020) and high-dimensional mediation analysis by  
22           Gao et al. (2019); while the superior methods for estimating the global mediation effect are high-  
23           dimensional linear mediation analysis by Zhou et al. (2021) and principal component mediation analysis  
24           by Huang and Pan (2016). We provide guidelines for epigenetic researchers on choosing the best method  
25           in practice and offer suggestions for future methodological development.

## 26 27           **Introduction**

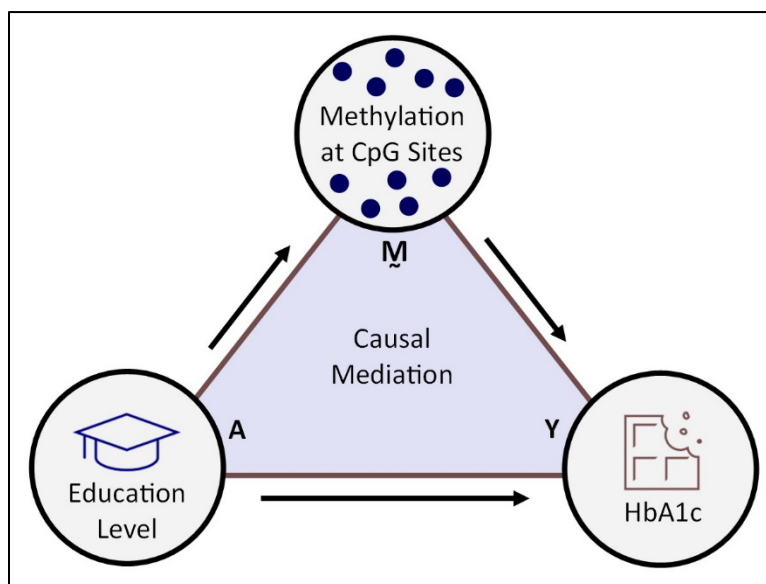
28           In this study, we review and evaluate the available methods for performing mediation analysis when the  
29           mediators are high-dimensional DNA methylation (DNAm) measurements. DNAm is an epigenomic  
30           mechanism describing when a methyl group binds to the DNA, which occurs predominantly at cytosine-  
31           guanine dinucleotides, called “CpG sites.” DNAm has an important role in regulating gene expression  
32           across the entire genome, and is particularly impactful at CpG sites in the promoter regions of genes,  
33           where it can inhibit the binding of enzymes needed for transcription<sup>1</sup>.

34           Recent advancements in technology have made it possible to collect DNAm data on a massive  
35           scale<sup>2</sup>. Indeed, microarray technologies have enabled the measurement of over 850,000 CpG sites  
36           simultaneously<sup>2</sup>, encouraging broad research on DNAm in the etiology of disease; and studies taking  
37           advantage of these tools have identified DNAm as a risk factor in obesity<sup>3,4</sup>, type II diabetes<sup>5</sup> and  
38           cardiovascular conditions<sup>6,7</sup>. At the same time, however, DNAm has also been linked to exposures such as  
39           diet<sup>8</sup>, smoking<sup>9</sup>, alcohol<sup>10</sup>, air pollution<sup>11</sup>, and socioeconomic status (SES)<sup>12,13</sup>, which has prompted  
40           research on whether the effects of these exposures on health outcomes could be transferred by changes in

41 DNAm. Effect transmission of this nature is called *mediation*, and it has become popular in epigenomic  
42 research to treat DNAm as a high-dimensional mediator between environmental exposures and human  
43 disease<sup>14</sup>.

44 As an example of such an analysis, our previous work<sup>15,16</sup> examined the association between low  
45 SES and glycated hemoglobin (HbA1c) in the Multi-Ethnic Study of Atherosclerosis (MESA), a United  
46 States population-based longitudinal study<sup>17</sup>. Indicators of SES, such as education level, are strong  
47 predictors of type II diabetes<sup>18</sup>, while HbA1c is an important risk factor of cardiovascular disease and a  
48 critical biomarker in type II diabetes diagnosis<sup>19-21</sup>. Since education level is also associated with  
49 DNAm<sup>12,13,22</sup>, and DNAm itself with HbA1c level<sup>23</sup>, we hypothesized that if low education results in  
50 greater HbA1c, part of that effect could be mediated by DNAm (Fig. 1). In the current study, we revisit  
51 this hypothesis for the purpose of illustration. Our sample from MESA has 963 individuals and includes  
52 DNAm measurements at 402,339 CpG sites, none of which we know for certain are related to education  
53 or HbA1c in advance.

54



55

56 **Fig. 1. Proposed causal mechanism in which the effect of low education on HbA1c is mediated by DNAm**

57

58 The standard statistical tool for addressing such a hypothesis is mediation analysis. Formally,  
59 mediation is when an exposure, say  $A$ , affects an outcome,  $Y$ , in part through its effect on a single  
60 mediating variable  $M$ . When  $M$  is a mediator of the  $A$  to  $Y$  association, the total effect of  $A$  on  $Y$  has two  
61 components: an *indirect effect*, from  $A$  affecting  $M$  and  $M$  affecting  $Y$ , and a *direct effect*, from  $A$  affecting  
62  $Y$  independently of  $M$ . In the “traditional mediation analysis” approach proposed by Baron and Kenny  
63 (1986), the associations from this mechanism could be measured by fitting a few regression models: one  
64 for the effect of  $A$  on  $M$  (the mediator model), one for the effects of  $A$  and  $M$  on  $Y$  (the outcome model),  
65 and sometimes a third model for the total effect of  $A$  on  $Y$ ,  $M$  ignored<sup>24-26</sup>. The more recently developed  
66 “causal mediation analysis,” based on the counterfactual approach<sup>27,28</sup>, has established conditions under  
67 which the parameters of these models can be interpreted as causal effects<sup>29</sup>. The causal approach is more  
68 flexible when  $Y$  or  $M$  are binary and when there is  $A$ - $M$  interaction in the outcome model<sup>30</sup>.

69

70

71 While standard examples of mediation consider only one exposure, one mediator, and one  
72 outcome<sup>31,32</sup>, there has been growing interest in methods for mediation that can handle many potential  
73 mediators at once. Epigenetic studies have felt this need especially, as DNAm is usually measured at  
74 several hundred thousand CpG sites with little prior knowledge of their importance. In settings such as  
75 this, a naïve strategy would be to evaluate the potential mediators one at a time, each with their own pair  
76 of models; but if the mediators are correlated this approach is inefficient, and the resulting estimates are  
77 potentially biased due to confounding from the excluded co-mediators<sup>31</sup>. Instead, so that we leverage  
78 these correlations rather than ignore them, the preferred approach is to assess the mediators jointly, in a  
79 single multivariable model. Although several methods for fitting such a model have been presented in the  
80 literature, none of them are widely used in analyzing DNAm data, a sign that epigenetic research is still  
81 catching up to recent developments in mediation analysis with high-dimensional mediators.

82 Our study aims to bridge this gap and guide researchers in epigenetics to use state of the art  
83 methods for mediation analysis with high-dimensional mediators. Despite the recent methodological  
84 developments, there are no clear-cut standards for which methods should be applied in which  
85 circumstances, making it difficult to select the best-suited method for an analysis in advance. While our  
86 prior research examined methods for large scale single-mediator hypotheses<sup>31</sup>, there is no such work for  
87 methods that can incorporate many potential mediators at once. Our study addresses this question first  
88 with an extensive simulation study, directly comparing the performance of seven different methods for  
89 mediation with high-dimensional mediators across a spectrum of settings. Along with metrics related to  
90 identification of key mediators and estimation of mediation effect, we include a computation time  
91 comparison to evaluate the scalability of the methods to large datasets. Next, to assess the utility of these  
92 methods on real data, we apply the same seven methods—plus two additional methods adapted from  
93 them—on the data from MESA to evaluate the mediating role of DNAm in the association between low  
94 education level and HbA1c. Our study is the first to address this critical gap in the epigenetic mediation  
95 literature, both by providing clarity on the methods available and by assessing their strengths and  
96 weaknesses under different settings. Moreover, although the analysis is centered around DNAm, the  
97 methods we deploy are not specific to epigenetics, and our results and guidelines should be similarly  
98 useful for researchers studying high-dimensional mediation problems in other fields. We include,  
99 supplementary to our study, an R package for implementing the methods, called “hdmed,” so that  
100 researchers have access to a centralized resource they can draw from in their own high-dimensional  
101 mediation analyses.

102

## 103 **Notations and General Framework**

104 Before proceeding, it will be useful to provide an overview of the relevant mediation model and  
105 to summarize the types of methods which have become available. To begin, suppose we have a dataset of  
106  $n$  individuals: an exposure  $A_i$ , a continuous outcome  $Y_i$ , and continuous mediators  $\mathbf{M}_i$  measured for the  $i^{\text{th}}$   
107 person,  $i$  varying from 1 to  $n$ . We write  $\mathbf{M}_i$  in bold to indicate its status as a vector—in this case, a set of  $p$   
108 mediators  $M_i^{(j)}$ ,  $j$  varying from 1 to  $p$ . Let  $\mathbf{C}_i$  be a vector of  $q$  covariates. When  $p$  is greater than 1, we can  
109 use the regression models

$$110 \quad E[Y_i|A_i, \mathbf{M}_i, \mathbf{C}_i] = \beta_a A_i + \boldsymbol{\beta}_m^T \mathbf{M}_i + \boldsymbol{\beta}_c^T \mathbf{C}_i \quad (1)$$

111 and

$$112 \quad E[\mathbf{M}_i|A_i, \mathbf{C}_i] = \boldsymbol{\alpha}_a A_i + \boldsymbol{\alpha}_c \mathbf{C}_i \quad (2)$$

113 to estimate the mediating role of  $\mathbf{M}_i$  in the causal pathway between the exposure and outcome<sup>33</sup>. Model  
114 (1) is the outcome model and model (2) is the mediator model. In model (1),  $\boldsymbol{\beta}_m$  is a  $p$ -vector in which the

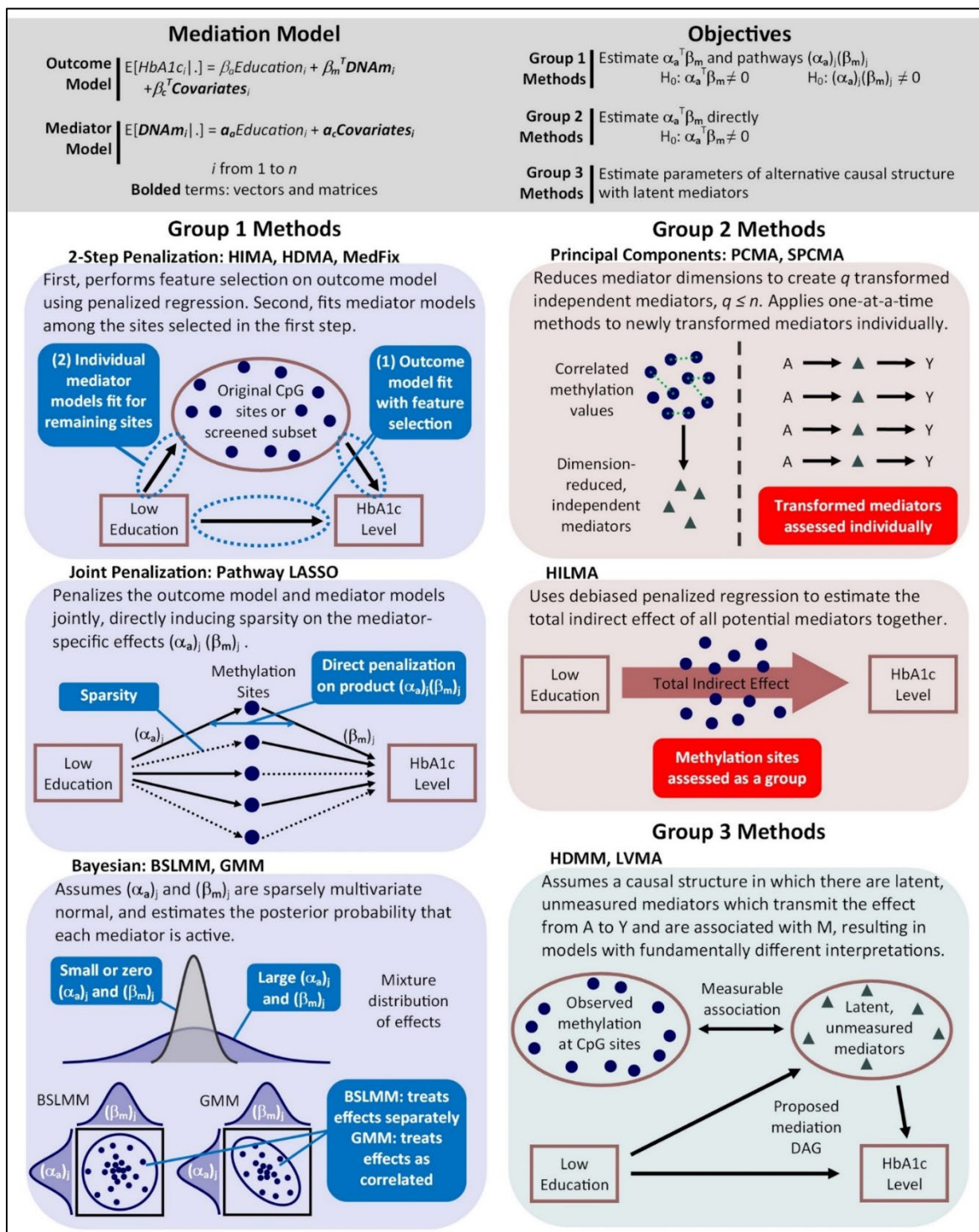
115  $j^{\text{th}}$  component,  $(\beta_m)_j$ , is the linear association of  $j^{\text{th}}$  mediator with  $Y_i$  adjusting for the other variables; while  
116  $\beta_a$  is the association between  $A_i$  and  $Y_i$  adjusting for mediators and covariates. In model (2),  $\alpha_a$  is a  $p$ -  
117 vector of the associations between the exposure and each mediator,  $(\alpha_a)_j$ ; and  $\alpha_c$  is a matrix with the  
118 mediator-covariate associations. Also note that in model (1), we have assumed there is no interaction  
119 between  $A_i$  and  $M_i$ , which is beyond the scope of our present study.

120 The parameters of these models underly the causal effects of interest. Under certain  
121 assumptions<sup>27,33</sup>, the direct effect of  $A_i$  on  $Y_i$  is  $\beta_a$ , the global indirect effect (or global mediation effect) of  
122  $A_i$  on  $Y_i$  through  $M_i$  is  $\alpha_a^T \beta_m$ , and the total effect of  $A_i$  on  $Y_i$  is  $\beta_a + \alpha_a^T \beta_m$ . Another quantity of interest is  
123 the proportion mediated, defined as the ratio of the global indirect effect to the total effect, which  
124 measures the degree to which the  $A_i$  to  $Y_i$  pathway is mediated by  $M_i$ . We may also seek to measure the  
125 product terms  $(\alpha_a)_j(\beta_m)_j$ , which measure the contribution of the  $j^{\text{th}}$  mediator to the global indirect effect,  
126 since summing these for  $j$  from 1 to  $p$  yields  $\alpha_a^T \beta_m$ . However, we emphasize that  $(\alpha_a)_j(\beta_m)_j$  cannot be  
127 interpreted as a causal effect through the  $j^{\text{th}}$  mediator on its own, since we have made no assumptions  
128 about the causal ordering of the mediators and can only formally treat them as a joint system. Instead, we  
129 call  $(\alpha_a)_j(\beta_m)_j$  the *mediation contribution*, and describe the  $j^{\text{th}}$  mediator as *active* if its contribution is not  
130 zero.

131 If the potential mediators are uncorrelated, conditional on the exposure and covariates, or if  $p$  is  
132 reasonably small relative to  $n$ , then it is trivial to fit the above models using linear regression. However, if  
133 the mediators are correlated and  $p$  is large, the estimates from model (1) may have extremely high  
134 variance; and if  $p$  is so large as to exceed  $n$ , the linear regression model cannot even be fitted. These  
135 concerns are relevant to us because DNAm measurements tend to be correlated, while the number of sites  
136 that we have measurements on exceeds the number of samples. Addressing these issues has been a focus  
137 of the mediation literature, with authors using penalized regression<sup>34–38</sup>, dimension reduction<sup>39–41</sup>,  
138 Bayesian inference<sup>15,42</sup>, and latent variables<sup>43</sup> to make the outcome model statistically tractable.

139 We provide a graphical depiction of eleven available methods in Fig. 2, dividing them into **three**  
140 different groups. Each method is described in greater detail in the Methods section and up to nine of them  
141 are included in the analysis. In the first group, we consider methods that fit the above pair of models  
142 explicitly, allowing one to estimate  $\alpha_a^T \beta_m$ , the global indirect effect, simply by summing the estimated  
143 mediation contributions. These include high-dimensional mediation analysis (HIMA) by Zhang et al.  
144 2016<sup>34</sup>, high-dimensional mediation analysis (HDMA) by Gao et al. 2019<sup>35</sup>, mediation analysis via fixed  
145 effect model (MedFix) by Zhang 2019<sup>36</sup>, pathway least absolute shrinkage operator (pathway LASSO) by  
146 Zhao and Luo 2022<sup>37</sup>, the Bayesian sparse linear mixed model (BSLMM) by Song et al. 2020<sup>15</sup>, and the  
147 Gaussian mixture model (GMM) by Song et al. 2021<sup>42</sup>. In the second group, we consider methods that  
148 can estimate  $\alpha_a^T \beta_m$  “directly”; in other words, without needing to fit the original pair of models explicitly.  
149 These have the drawback of being unable to identify specific active mediators because they do not  
150 provide estimates of the mediation contributions. They include principal component mediation analysis  
151 (PCMA) by Huang and Pan 2016<sup>39</sup>, sparse principal component mediation analysis (SPCMA) by Zhao et  
152 al. 2020<sup>40</sup>, and high-dimensional linear mediation analysis (HILMA) by Zhou et al. 2021<sup>38</sup>. Last, in the  
153 third group, we consider methods that make no attempt to estimate the mediation effects as originally  
154 proposed, but instead reconceptualize the mediation framework with newly-defined parameters based on  
155 latent variables. This group includes the methods high-dimensional multivariate mediation analysis  
156 (HDMM) by Chén et al. 2018<sup>41</sup> and latent variable mediation analysis (LVMA) by Derkach et al. 2021<sup>43</sup>.  
157 Within this comparative structure, we evaluate methods from all three groups, identifying their strengths  
158 and weaknesses across a wide range of simulation settings and analysis of DNAm data from MESA.

159  
160  
161



162 **Fig. 2. Methods for mediation analysis with high-dimensional DNAm data.** Figure describes eleven methods for  
 163 mediation analysis that can be applied to high-dimensional DNA methylation data, each of which is described in  
 164 greater detail in the Methods section. Seven of these methods are included in the simulation study and nine in the

165 observed DNAm data analysis with MESA. Group 1 methods fit the outcome model explicitly using penalized  
166 regression or Bayesian regression; Group 2 methods obtain the global mediation effect without fitting the original  
167 outcome model explicitly; and Group 3 methods measure mediation through latent variables.

168

## 169 Results

### 170 Simulation Results

171 We begin by comparing the performance of the methods using simulations, where we know and  
172 can control the true values of the parameters. On simulated data with 2,000 (potential) mediators and  
173 either 1,000 or 2,500 observations, we consider (1) a baseline setting, where the mediators are moderately  
174 correlated and their signals are sparse; (2) a high-correlation setting, where the correlations between  
175 mediators are enhanced compared to (1); and (3) a non-sparse setting, where every mediator has at least  
176 some mediation signal but some of the signals are systematically larger. In Settings (1) and (2), 60  
177 random mediators have  $(\alpha_a)_j$  only sampled from a Normal(0,1), 60 have  $(\beta_m)_j$  only sampled from a  
178 Normal(0,1), and 20 have both, with the remaining entries of  $\alpha_a$  and  $\beta_m$  fixed at zero. In Setting (3), we  
179 use a similar scheme, but sample the previously zero  $(\alpha_a)_j$  and  $(\beta_m)_j$  from a Normal(0,0.2<sup>2</sup>). Our  
180 simulations also vary the strength of the signals within each of these settings by changing the proportion  
181 of variance that is explained by the associations. We do so by changing PVE<sub>A</sub>, the proportion of variance  
182 in each mediator that can be explained by *A*, among those mediators that are affected by *A*; PVE<sub>IE</sub>, the  
183 proportion of variance of *Y* that is explained by the total mediation effect; and PVE<sub>DE</sub>, the proportion of  
184 variance of *Y* that is explained by the direct effect of *A* on *Y*. Results for varying PVE<sub>IE</sub> are presented here  
185 while results for varying PVE<sub>DE</sub> and PVE<sub>A</sub> are included in the supplement (Supplementary Figs 1-4). In  
186 addition to the high-dimensional mediation methods, we include a one-at-a-time method<sup>44</sup> in which the  
187 mediators are assessed individually using linear regression. We evaluate the methods by their true  
188 positive rate (TPR) for detecting active mediators, their mean squared error (MSE) for estimating the  
189 contributions of active and inactive mediators, and their percent relative bias for estimating the global  
190 indirect effect. See Methods for more details.

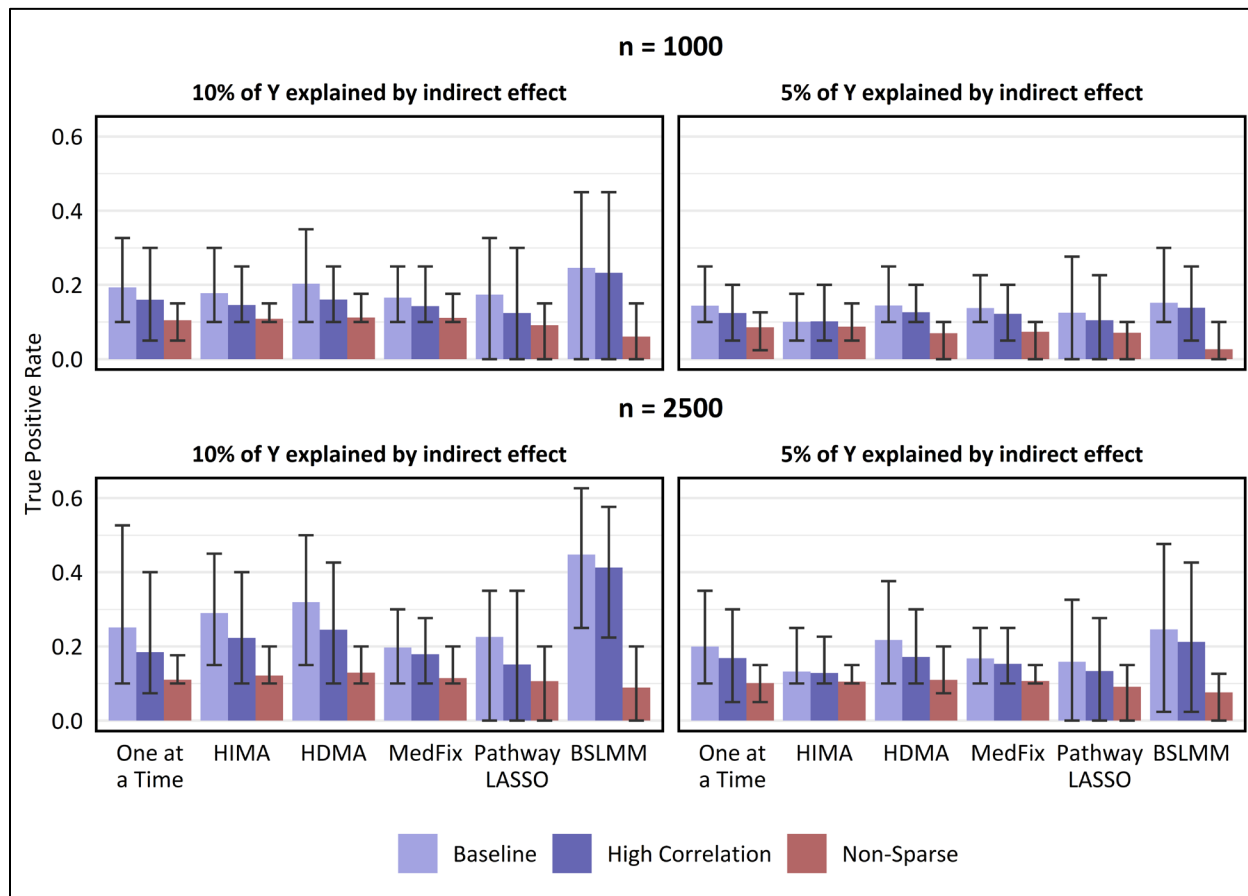
191

#### 192 True positive rate

193 Fig. 3 compares the TPR detecting active mediators of the Group 1 methods and the one-at-a-time  
194 method. The value shown is the mean TPR over 100 simulated datasets and a 95% empirical confidence  
195 interval (CI). On each dataset and for each method, thresholding was used to keep the false discovery rate  
196 (FDR) below 10%. For the non-sparse setting, we show the TPR for detecting mediators whose  $(\alpha_a)_j$  and  
197  $(\beta_m)_j$  were both sampled from Normal(0,1) rather than Normal(0,0.2<sup>2</sup>). We include the Group 1 methods  
198 HIMA, HDMA, MedFix, pathway LASSO, and BSLMM. We focus on TPR but not false positive rate  
199 (FPR) because the FDR correction was highly conservative, the mean FPR ranging from 0 to 5.1x10<sup>-4</sup>  
200 across all settings and methods.

201 For a sample size of 2,500 and a PVE<sub>IE</sub> of 0.10, the most powerful method in the baseline setting  
202 was BSLMM (mean TPR: 0.45; CI: 0.25 - 0.63), whose average TPR was 40% higher than that of the  
203 second-best method, HDMA. BLSMM also performed best when PVE<sub>IE</sub> was 0.05 (mean TPR: 0.25; CI:  
204 0.02 - 0.48), but to a lesser degree, outperforming HDMA by only 13%. BSLMM remained the best  
205 method, and HDMA the second best, no matter the signal strength or the degree of correlations, but  
206 performed poorly when the signals were non-sparse. In the setting with 1,000 observations, PVE<sub>IE</sub> set to

207 0.05, and non-sparse signals, the best-performing method was HIMA (mean TPR: 0.09; CI: 0.05 - 0.10),  
 208 its average TPR 3.3 times higher than that of BSLMM, which performed worst.  
 209

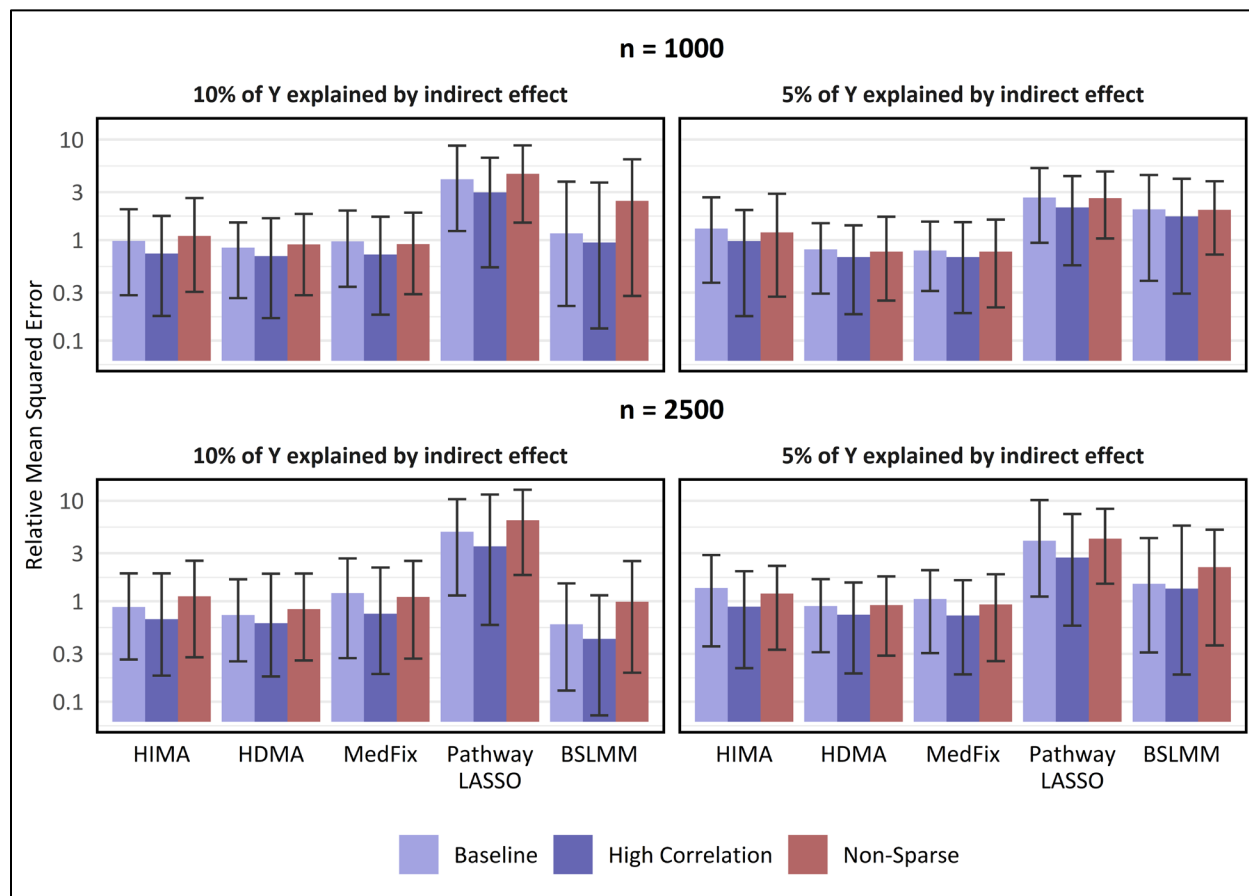


210 **Fig. 3. True positive rate for detecting mediation signals at a false discovery rate of 10%.** Value shown is the  
 211 mean TPR across 100 simulated data replicates, with intervals representing the inner 95% range. In the baseline and  
 212 high-correlation-among-mediators settings, TPR is for distinguishing mediators which contribute to the global  
 213 mediation effect from those which do not, whereas in the non-sparse setting, TPR is for distinguishing mediators  
 214 whose contributions were sampled from a high-variance distribution from those whose contributions were sampled  
 215 from a low-variance distribution. False discovery proportion was capped below 10% by a proper choice of the p-  
 216 value threshold (one-at-a-time, HIMA, HDMA, MedFix), posterior inclusion probability threshold (BSLMM), or  
 217 method tuning parameter (pathway LASSO).  
 218

### 219 Estimation of contributions of active mediators

220 Next, we assess the MSE of the methods for estimating mediation contributions, relative to the one-at-a-  
 221 time approach. In Fig. 4, we show the relative MSE (rMSE) for estimating mediation contributions  
 222 among the mediators that were either active (in the baseline and high-correlation settings) or had ( $\alpha_m$ ) or  
 223 ( $\beta_m$ ) sampled from the larger-variance distribution (in the non-sparse setting). In the baseline setting with  
 224 2,500 observations, the best-performing method when the mediation signal was strong was BSLMM,  
 225 whose mean rMSE of 0.59 (CI: 0.13 - 1.51) was 24% lower than that of HDMA, the second-best method.  
 226 However, when the PVE<sub>IE</sub> was reduced to 0.05 or the sample size reduced to 1,000, the best-performing  
 227 method was either HDMA or MedFix, with MedFix (mean rMSE: 0.79; CI: 0.31 - 1.53) performing 61%

228 better than BSLMM after reducing both. Similar trends were observed for the high-correlation and non-  
 229 sparse settings.  
 230



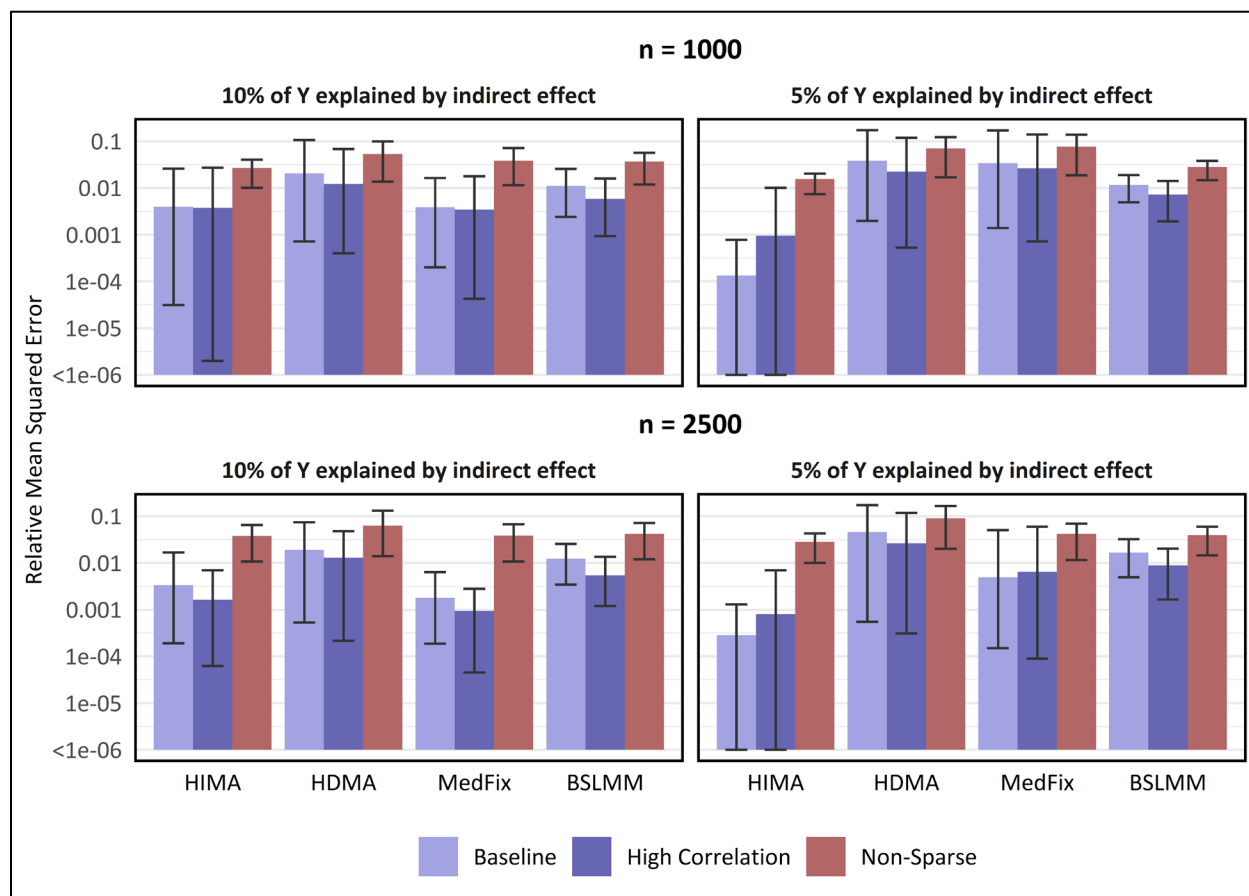
231 **Fig. 4. MSE in estimating mediation contributions of active mediators, relative to one-at-a-time method.** Y-  
 232 axis is on a log<sub>10</sub> scale. Value shown is the mean of the relative mean-squared error for estimating mediation  
 233 contributions among active mediators (relative to the one-at-a-time approach) across 100 simulated data replicates,  
 234 with intervals representing the inner 95% range. For baseline and high-correlation-between-mediators settings,  
 235 active mediators which contribute to the global mediation effect, whereas in the non-sparse setting, active mediators  
 236 are those whose contributions were sampled from a distribution with large variance instead of small.  
 237

### 238 Estimation of contributions of inactive mediators

239 Figure 5 shows the rMSE among the mediators that either were not active (in the baseline and  
 240 high-correlation settings) or had  $(\alpha_i)_j$  or  $(\beta_m)_j$  sampled from the smaller-variance distribution (in the non-  
 241 sparse setting). We exclude pathway LASSO from Fig. 4 because for the baseline and high-correlation  
 242 settings it had rMSEs of exactly zero. The reason for this is that pathway LASSO tended to be highly  
 243 conservative and successfully assigned inactive mediators to have no effect. As for the other methods, in  
 244 the baseline setting with 2,500 samples, MedFix performed the best when PVE<sub>IE</sub> was 0.10, with a mean  
 245 rMSE of  $1.8 \times 10^{-3}$  (CI:  $1.9 \times 10^{-4} - 6.4 \times 10^{-3}$ ), which was 46% lower than the mean rMSE for the second-best  
 246 method, HIMA. In contrast, HIMA was the best-performing method when signal was weakened to a  
 247 PVE<sub>IE</sub> of 0.05, attaining a mean rMSE of  $2.8 \times 10^{-4}$  (CI:  $0.0 - 1.3 \times 10^{-3}$ ), which was 94% lower than that of  
 248 the second-best, MedFix. Results were similar when the correlations between mediators were heightened  
 249 and when the sample size was reduced. In the settings where mediation signals were non-sparse, the best-



250 performing method was always HIMA, which had a mean rMSE of  $3.7 \times 10^{-2}$  (CI:  $1.1 \times 10^{-2} - 6.5 \times 10^{-2}$ )  
 251 when  $PVE_{IE}$  was 0.10 and there were 2,500 observations, 2% lower than that of MedFix.  
 252

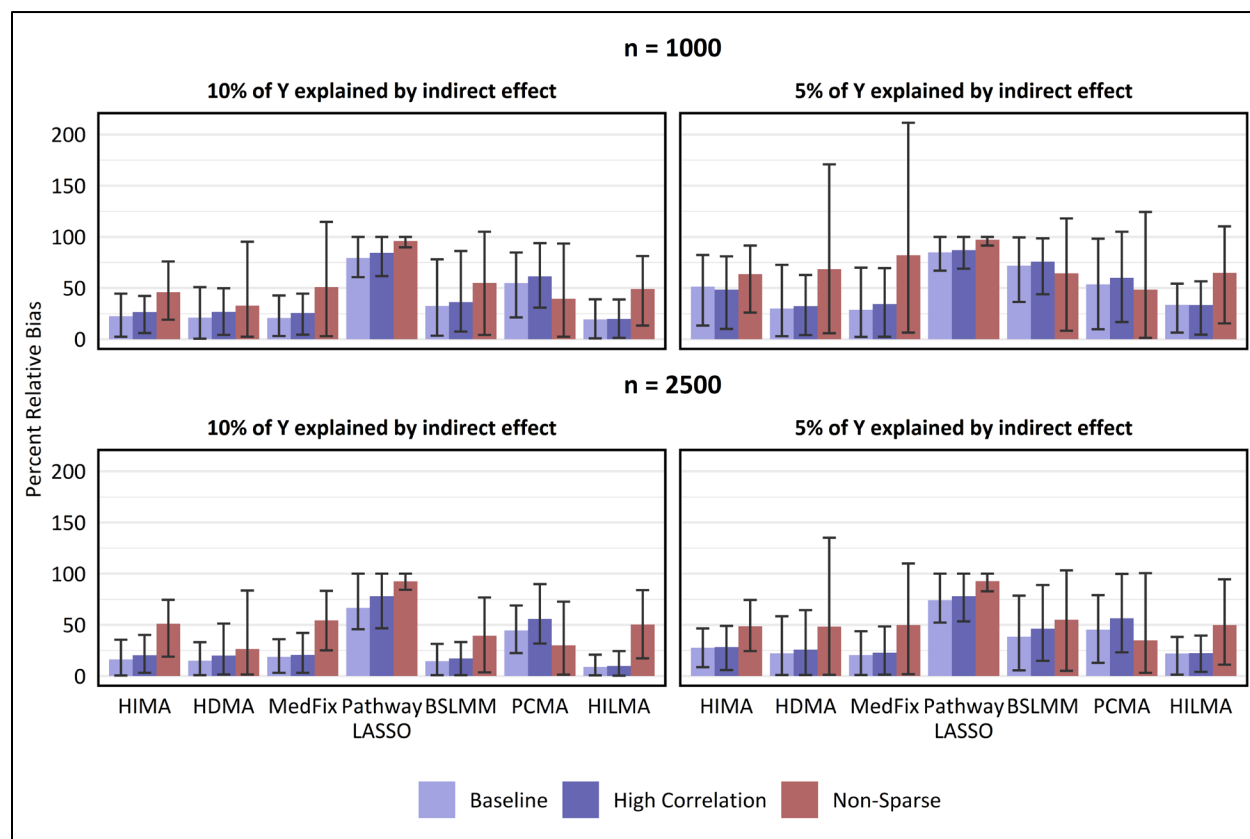


253 **Fig. 5. MSE in estimating mediation contributions of inactive mediators, relative to one-at-a-time method.**  
 254 Y-axis is on a log<sub>10</sub> scale. Value shown is the mean of the relative mean-squared error for estimating mediation  
 255 contributions among inactive mediators (relative to the one-at-a-time approach) across 100 simulated data replicates,  
 256 with intervals representing the inner 95% range. For baseline and high-correlation-between-mediators settings,  
 257 inactive mediators are those which do not contribute to the global mediation effect, whereas in the non-sparse  
 258 setting, inactive mediators are those whose contributions were sampled from a distribution with small variance  
 259 instead of large.

260  
 261 **Estimation of global indirect effect**

262 Lastly in Fig. 6, we show the percent relative bias for estimating  $\alpha_a^T \beta_m$ , the global indirect effect.  
 263 We use the same methods as in Figures 3 to 5 along with the Group 2 methods PCMA and HILMA,  
 264 which obtain an estimate of the global indirect effect without needing to directly fit the original mediation  
 265 model. (The Group 2 method SPCMA is excluded for computational reasons.) In the baseline setting with  
 266 2,500 samples, the best performer when  $PVE_{IE}$  was 0.10 was HILMA, whose mean relative bias of 9%  
 267 (CI: 0.6% - 20.8%) was 40% lower than that of HDMA, the second-best. Next, when the PVE was  
 268 reduced to 0.05, the best-performing method was MedFix (mean relative bias: 20.5%; CI: 1.0% - 43.8%),  
 269 which outperformed HILMA by only 7%. We observed similar results for a sample size of 1,000 and  
 270 high-correlations. In the non-sparse settings, where the biases tended to be much higher, the best  
 271 performing methods were either PCMA or HDMA.

272  
273



274 **Fig. 6. Percent relative bias in estimated global indirect effect.** Value shown is the mean of the percentage  
 275 relative bias in estimating the global mediation effect across 100 simulated data replicates, with intervals  
 276 representing the inner 95% range.

277

## 278 Scalability

279 We evaluated the scalability of the methods by running them 30 times on a common computing  
 280 platform, and recording their run time (Table 1). This was done in both a small data setting ( $n = 100, p =$   
 281  $200$ ) and a big data setting ( $n = 1,000, p = 1,000$ ). On the larger dataset, the methods MedFix, HDMA,  
 282 and PCMA posed insignificant computational burden; whereas BSLMM took an average of 40.1 minutes  
 283 per run (assuming 30,000 posterior samples), HILMA an average of 40.9 minutes per run, pathway  
 284 LASSO an average of 192.6 minutes per run, and SPCMA an average of 842.5 minutes per run (assuming  
 285 100 principal components). Run times were substantially lower in the smaller dataset, the slowest method,  
 286 pathway LASSO, only taking an average of 18.71 minutes. The memory consumption of the methods is  
 287 included in Supplementary Table 1.

288

289

290

291

292

293

294

295  
296  
297

**Table 1. Computation time comparison for high-dimensional mediation analysis methods**

Method	$n = 100, p = 200$		$n = 1,000, p = 2,000$	
	Mean	Interquartile Range	Mean	Interquartile Range
BSLMM	39.17s	(38.84s - 39.54s)	40.14m	(39.74m - 40.34m)
HDMA	1.40s	(1.37s - 1.40s)	29.76s	(29.55s - 29.92s)
HDMM	24.85s	(24.80s - 24.89s)	12.36m	(12.33m - 12.37m)
HILMA	24.42s	(24.13s - 24.63s)	40.85m	(38.22m - 40.65m)
HIMA	0.25s	(0.25s - 0.25s)	3.55s	(3.47s - 3.62s)
MEDFIX	0.61s	(0.60s - 0.61s)	7.33s	(7.22s - 7.42s)
PCMA	2.77s	(2.74s - 2.79s)	58.97s	(58.08s - 59.35s)
PLASSO	18.71m	(18.19m - 19.23m)	192.62m	(188.10m - 195.83m)
SPCMA	16.05m	(15.94m - 16.04m)	842.54m	(827.26m - 855.21m)

298 Methods were run 30 times each on a single core of an Intel(R) Xeon(R) Gold 6242R CPU @ 3.10GHz processor.

299

## 300 DNAm data analysis results from MESA

301 For our real data analysis, we applied the methods on a dataset with high-dimensional epigenetic  
302 mediators. Our exposure of interest was low SES—measured by educational attainment below a four-year  
303 degree—while our outcome variable was HbA1c level and our potential mediators were DNAm  
304 measurements at 402,339 CpG sites. Since the methods are incapable of handling so many CpG sites at  
305 once, we reduced our scope to only include the 2,000 sites with the strongest association with low SES.  
306 This was based on a linear mixed-model adjusting for age, sex, race, and the estimated proportions of  
307 residual non-monocytes as fixed effects and methylation chip and position as random effects. Our final  
308 dataset contained these 2,000 CpG sites and 963 samples. HbA1c, DNAm, and all other continuous  
309 variables were standardized prior to analysis.

310

### 311 Identification of noteworthy CpG sites

312 We identified CpG sites that potentially mediated the relationship between low SES and HbA1c  
313 using the Group 1 methods HIMA, HDMA, MedFix, pathway LASSO, and BSLMM. In HIMA, HDMA,  
314 MedFix, and pathway LASSO, which involve feature selection, we describe a CpG site to be “active” if  
315 its estimated mediation contribution is not zero; whereas in BSLMM, we do so if the estimated posterior  
316 inclusion probability is not zero (see Methods). We also included a one-at-a-time method in which the  
317 CpG sites were assessed individually with linear mixed models, identifying active mediators with the  
318 joint significance test<sup>44</sup>. Out of 2,000 CpG sites, HIMA found 3 sites to be noteworthy, HDMA found 11,  
319 MedFix found 3, pathway LASSO found 141, and BSLMM found 3, amounting to 144 unique CpG sites  
320 in total. The one-at-a-time method identified zero CpG sites as noteworthy at an FDR threshold of 10%.  
321 Eleven CpG sites were identified as noteworthy by at least two of the methods (Table 2). Among these  
322 11, the estimated mediation contributions were similar across methods in direction and size except for  
323 BSLMM, for which the estimates were an order of magnitude smaller than the others but in the same  
324 direction.

325

326

327

328

329

330  
331  
332  
333  
334

**Table 2. Estimated contributions of noteworthy CpG sites on the mediation pathway between low education and HbA1c**

CpG Name	Chromosome	Nearby Gene(s)	USCS RefGene Group	Univariate (0 sites identified)	HIMA (3 sites identified)	HDMA (11 sites identified)	MedFix (3 sites identified)	Pathway LASSO (141 sites identified)	BSLMM (3 sites identified)
cg10508317	17	<i>SOCS3</i>	Body	$3.48 \times 10^{-2}$	$1.59 \times 10^{-2*}$	$3.56 \times 10^{-2*}$	$2.90 \times 10^{-2*}$	$2.35 \times 10^{-2*}$	$0.25 \times 10^{-2}$
cg01288337	14	<i>RIN3</i>	Body	$3.35 \times 10^{-2}$	$1.47 \times 10^{-2*}$	$2.82 \times 10^{-2*}$	$2.70 \times 10^{-2*}$	$4.43 \times 10^{-2*}$	$0.21 \times 10^{-2}$
cg10244976	16	<i>LMF1</i>	Body	$3.00 \times 10^{-2}$	0	$2.78 \times 10^{-2*}$	0	$2.23 \times 10^{-2*}$	$0.19 \times 10^{-2}$
cg07516252	14	<i>REC8</i>	TSS200	$2.72 \times 10^{-2}$	0	$2.24 \times 10^{-2*}$	0	$2.26 \times 10^{-2*}$	$0.26 \times 10^{-2}$
cg07571519	10	<i>C10orf105</i> ; <i>CDH23</i>	3'UTR; Body	$2.53 \times 10^{-2}$	$0.33 \times 10^{-2*}$	$3.67 \times 10^{-2*}$	$1.47 \times 10^{-2*}$	$2.81 \times 10^{-2*}$	$0.21 \times 10^{-2}$
cg23079012	2	<i>LINC00299</i>	Body	$2.27 \times 10^{-2}$	0	$1.99 \times 10^{-2*}$	0	$1.98 \times 10^{-2*}$	$0.29 \times 10^{-2}$
cg01587454	8	<i>DCAF4L2</i>	1stExon	$1.77 \times 10^{-2}$	0	$2.10 \times 10^{-2*}$	0	$1.99 \times 10^{-2*}$	$0.38 \times 10^{-2}$
cg27527503	4	<i>HADH</i>	TSS1500	$1.75 \times 10^{-2}$	0	$1.86 \times 10^{-2*}$	0	$1.27 \times 10^{-2*}$	$0.23 \times 10^{-2}$
cg25891647	11	<i>GRAMD1B</i>	Body	$-1.27 \times 10^{-2}$	0	$-3.42 \times 10^{-2*}$	0	$-3.02 \times 10^{-2*}$	$-0.33 \times 10^{-2}$
cg08473752	17	<i>NLK</i>	Body	$-0.70 \times 10^{-2}$	0	$-2.34 \times 10^{-2*}$	0	$-2.32 \times 10^{-2*}$	$-0.22 \times 10^{-2}$
cg12644059	15	<i>BLM</i>	N/A <sup>1</sup>	$-0.03 \times 10^{-2}$	0	$-2.31 \times 10^{-2*}$	0	$-1.84 \times 10^{-2*}$	$-0.22 \times 10^{-2}$

335 \*Selected as noteworthy by given method

336 <sup>1</sup>CpG site cg12644059 is 3.240kb from the final base pair of the BLM gene

337 Table includes all CpG sites that were selected as having a noteworthy mediation contribution by at least two of the  
338 implemented methods out of 2,000 CpG sites in total. Criteria for CpG identification varied by method. All  
339 estimates are adjusted for age, sex, race, and the estimated proportions of residual non-monocytes as fixed effects,  
340 along with methylation chip and position as random effects to address potential batch effects. Note that for HIMA,  
341 HDMA, MedFix, and pathway LASSO, which fit high-dimensional regression models, we used additional pre-  
342 screening to reduce the number of mediators in advance to only  $n/\log(n) \approx 141$  CpG sites, which is the approach  
343 recommended by the HIMA and HDMA authors and helps with statistical and computational efficiency (see  
344 Methods). Pathway LASSO selected all of these 141.

345

346 Some of these CpG sites were on or nearby genes that are potentially related HbA1c. Site  
347 cg10508317 is in the body of the *SOCS3* gene, for which a rich body of literature has established links  
348 between overexpression and insulin resistance<sup>45</sup>. The same site has also been identified in MESA as a  
349 mediator between adult SES and BMI<sup>46</sup> and adult SES and HbA1c<sup>31</sup> based on previous one-at-a-time  
350 analyses. Site cg01288337, in the body of the *RIN3* gene, has been identified in MESA as a potential  
351 mediator between adult SES and HbA1c based on one-at-a-time analysis as well<sup>31</sup>. The *RIN3* gene itself is  
352 proximal to the *SLC24A4* gene, both of which have been linked to brain glucose metabolism in human  
353 population studies<sup>47</sup>. In addition, site cg27527503 is in the promoter region of the *HADH* gene, which is  
354 differentially expressed with respect to diabetes status<sup>48</sup> and is a primary driver of hyperinsulinism<sup>49</sup> and  
355 hyperinsulinaemic hypoglycemia (low blood sugar due to excess insulin)<sup>50</sup>. A Venn diagram of genes  
356 identified by the methods is included in Supplementary Fig. 5, and results for every noteworthy CpG site  
357 are listed in Supplement File 1.

358

### 359 Global mediation through DNAm

360 Next, we estimated the direct effect of low education on HbA1c, the global indirect effect of low  
361 education on HbA1c through DNAm, and the total effect of low education on HbA1c using the Group 1

362 methods HIMA, HDMA, MedFix, pathway LASSO, and BSLMM, as well as the Group 2 methods  
363 PCMA, SPCMA, and HILMA (Table 3). Results across methods varied considerably, with the estimated  
364 global indirect effect ranging from 0.03 in HILMA to 0.17 in SPCMA. The estimated total effect ranged  
365 from 0.02 (HILMA) to 0.198 (HIMA, HDMA, and MedFix). While HILMA appeared to be an outlier,  
366 some of the other methods were consistent, with HDMA, BSLMM, P- LASSO, PCMA, and SPCMA all  
367 estimating the global indirect effect to be close to 0.15. The variability in the estimated indirect effect and  
368 estimated total effect led to variability in the proportion mediated as well, from 17.1% in HIMA to 100%  
369 in HILMA.

370

371 **Table 3. Estimated effects in the mediation mechanism from low education to DNAm to HbA1c**

Method	Estimated Global indirect Effect	Estimated Direct Effect	Estimated Total Effect	Estimated Proportion Mediated
HIMA	0.03	0.16	0.20	0.17
HDMA	0.13	0.07	0.20	0.65
MedFix	0.07	0.13	0.20	0.36
BSLMM	0.14	0.05	0.18	1.00
Pathway LASSO	0.13	0.05	0.18	0.74
PCMA	0.15	0.02	0.17	0.91
SPCMA	0.17	0.00	0.17	1.00
HILMA	0.03	0.00	0.03	1.00

372 All estimates are adjusted for age, sex, race, and the estimated proportions of residual non-monocytes as fixed  
373 effects, along with methylation chip and position as random effects to address potential batch effects. We provide  
374 only point estimates, not interval estimates, because some of the methods are either not capable of producing  
375 interval estimates or do not provide the code for producing them in their software. For HIMA, HDMA, and MedFix,  
376 which as coded do not directly provide estimates of the direct effect, we first estimate the total effect by fitting the  
377 outcome model with the CpG sites omitted, then estimate the direct effect by subtracting the indirect effect from the  
378 total effect. Note also that, for HIMA, HDMA, MedFix, and pathway LASSO, we used additional screening to  
379 reduce the number of mediators in advance for the sake of statistical and computational efficiency, so only  $n/\log(n)$   
380  $\approx 141$  CpG sites were seen by the multivariate model rather than 2,000 (this approach is recommended by the HIMA  
381 and HDMA authors).

382

### 383 **Additional Findings**

384 In addition to estimating the global indirect effect, method SPCMA is also able to identify  
385 potentially-mediating CpG sites in groups. It does so by linearly combining the mediators using sparse  
386 principal component-defined weights, then evaluating the resulting principal components as mediators  
387 themselves<sup>40</sup>. However, out of 100 computed principal components, only three of them had significant  
388 mediation contributions after 10% FDR correction, the first representing a linear combination of 762 CpG  
389 sites, the second a combination of 782 sites, and the third a combination of 797 sites. Since the  
390 transformed mediators are functions of so many CpG sites at once, one cannot make claims about which  
391 particular CpG sites are active mediators, but the method still provides insight to whether there is  
392 statistical mediation at all.

393 We finish our analysis by deploying HDMM, a method from Group 3. Unlike the methods in  
394 Groups 1 and 2, HDMM cannot be used to estimate the global indirect effect from the proposed mediation  
395 structure, nor to estimate the mediation contributions of specific CpG sites. Rather, HDMM uses a  
396 likelihood-based approach to compute “directions of mediation”, which are weights that can be used to  
397 linearly combine the observed mediators into unobserved, latent mediators that replace the observed  
398 mediators in the mediation models (similar to PCMA). The estimated effect of the first latent mediator on

399 average HbA1c was 0.13, the estimated total effect 0.71, and the proportion mediated 0.715. The three  
400 CpG sites with the largest directions of mediation were cg01288337 (0.36) on the *RIN3* gene,  
401 cg16162970 (-0.22) near the *PACS2* gene, and cg25891647 (-0.21) on the *GRAMD1B* gene; the first and  
402 last of which were among the 11 CpG sites identified by other methods in Table 2. Although the size and  
403 direction of these estimates are not interpretable, they offer evidence that these CpG sites are potentially  
404 involved in mediation.

405

## 406 Discussion

407 In this study, we reviewed and evaluated statistical methods for performing mediation analysis with high-  
408 dimensional DNAm data, so that researchers in epigenetics have the information they need to choose the  
409 most appropriate method for their data sample, subject matter, and research objectives. In extensive  
410 simulations, we found that the most powerful method for identifying active mediators was generally  
411 BSLMM, with HDMA close behind; though the former performed poorly in settings where the mediation  
412 signals were non-sparse. No method was uniformly better than the others at estimating the mediation  
413 contributions, though pathway LASSO was always the weakest. For estimating the global indirect effect,  
414 the best-performing method was HILMA in sparse mediation settings and PCMA or HDMA in non-  
415 sparse settings. Our scalability comparison revealed that HIMA, HDMA, MedFix, and PCMA were easily  
416 scalable to large datasets (e.g.,  $n = 1,000$  and  $p = 2,000$ ), whereas SPCMA and pathway LASSO were  
417 extremely computationally costly.

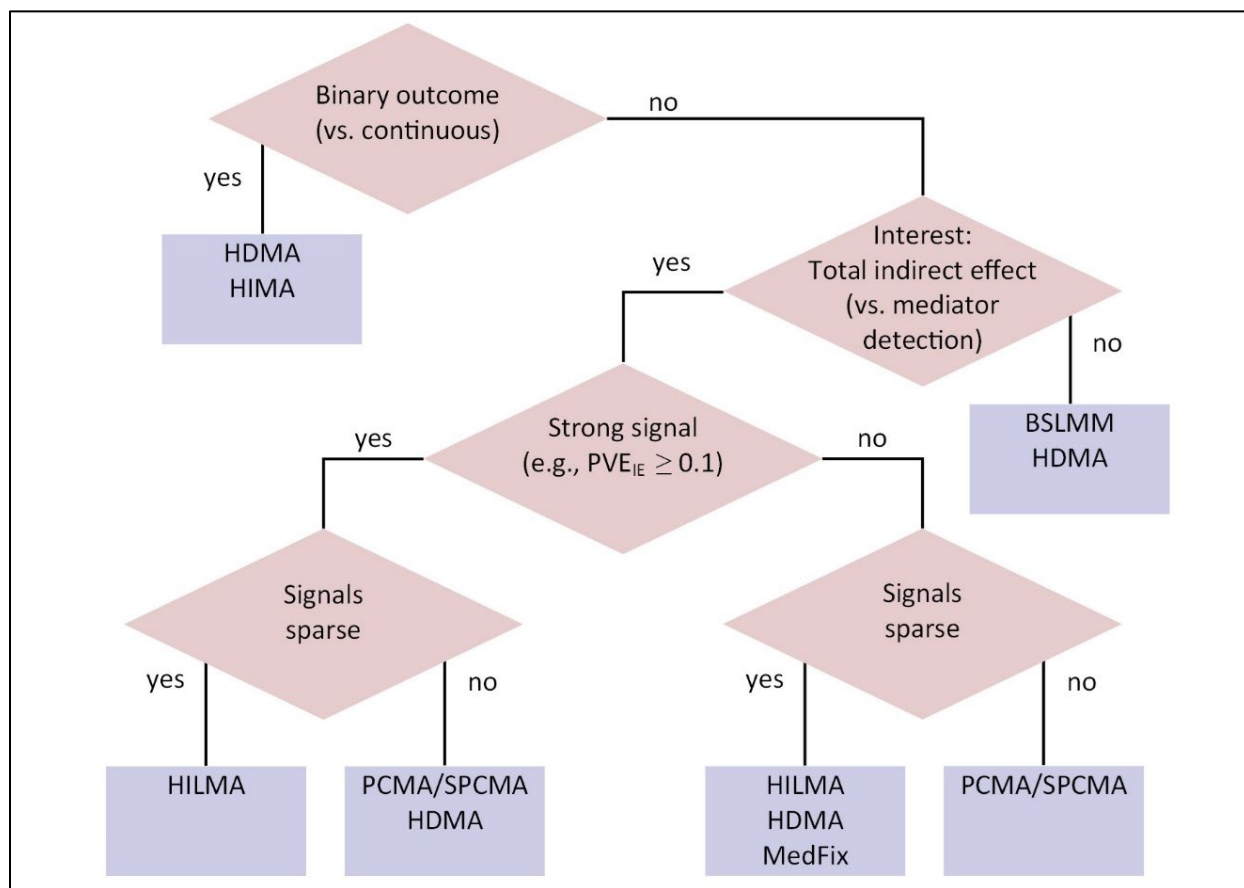
418 On DNAm data from MESA, 11 CpG sites were selected by at least two of the methods as  
419 mediators between low SES and HbA1c level. Of the many genes related to these sites, *SOCS3*, *RIN3*,  
420 and *HADH* have the strongest potential biological connections to HbA1c<sup>45,47,48,50-52</sup>, which contributes to  
421 the already rich literature on DNAm as a mediator between the exposome and health outcomes.  
422 Moreover, the methods generally produced similar estimates of the mediation contributions, with the  
423 exception of BSLMM. It is possible that since BSLMM is non-sparse, the estimated mediation  
424 contributions end up severely shrunken compared to the methods which directly select features.

425 Estimates of the global indirect effect were highly variable. Part of this can be explained by the  
426 fact that HDMA, MedFix, HIMA, and pathway LASSO are sparse models that can set mediation  
427 contributions to be exactly zero, resulting in a rigid and unstable estimation of the global indirect effect.  
428 The method HILMA, which is built specifically for estimating the global indirect effect and direct effect,  
429 produced estimates that were sharply different than the other methods, possibly because our simulations  
430 indicated that it struggled in non-sparse mediation settings.

431 In practice, the optimal method for mediation analysis with high-dimensional mediators will  
432 depend both on the data and the objective. If the goal is to identify specific CpG sites that are involved in  
433 mediation, one preferred method may be HDMA, which performed well at detecting active mediators in  
434 our simulations and was not overly conservative when applied to the observed data. If one's focus is the  
435 global indirect effect, our simulations suggested that the optimal method is HILMA; but considering the  
436 variability we observed in our DNAm analysis, it may be worthwhile to apply BSLMM and HDMA as  
437 well to ensure the results are robust. If the results of multiple methods disagree substantially, it may be  
438 difficult to say with confidence which is closest to the truth, and the estimates should be interpreted with  
439 caution. Next, if there is interest in latent, unmeasured mediators, either HDMM or LVMA is worth  
440 attempting, though HDMM is computationally simpler. A detailed decision tree for selecting the optimal  
441 method is included in Fig. 7.

442

443  
444



445 **Fig. 7. Decision tree for selecting a high-dimensional mediation analysis.**

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

Some strengths of our study include its broad coverage of the available methods, the breadth of its simulation settings, and the comprehensive set of evaluation criteria. Our analysis of real DNAm data is especially essential because it elucidates the potential limitations of using these methods in practice, as it is impossible to incorporate the full complexity of real data sources into contrived simulation settings. However, our study also has weaknesses. First, since DNAm measurements and HbA1c data were collected concurrently, and represent only single time points, we cannot interpret the parameters we have estimated as causal effects. Nor can we interpret the mediation contributions estimated in Table (2) as causal, since DNAm was correlated across CpG sites and we have made no assumptions about their causal ordering. Moreover, although it would be optimal to address our research question longitudinally, with measurements at multiple time points, there is a dearth of mediation analysis methods which can handle that type of data, and longitudinal mediation analysis with high-dimensional mediators should be a focus of future methodological development. Second, we limited our analysis to the situation that  $Y$  and  $M$  are continuous, that  $M$  and  $A$  do not interact, and that only one  $A$  is of interest. However, we note that the methods HIMA and HDMA can also be applied to identify active mediators when  $Y$  is binary, while PCMA can be applied to infer the global indirect effect when there is  $A$ - $M$  interaction in the outcome model. MedFix, along with the simultaneously-proposed MedMix (mediation analysis with mixed effect model by Zhang (2021)) can be applied when both the exposures and mediators are high-dimensional,

464 while Huang and Vanderweele (2014) proposed a variance component test of the global indirect effect  
465 when only  $A$  is high-dimensional<sup>53</sup>. As the landscape of methods for high-dimensional mediation analysis  
466 continues to expand, future review studies should consider exploring additional mediation settings (in  
467 presence of non-linearity, interaction) for which statistical methods are continuing to become available.  
468

## 469 **Methods**

### 470 **Mediation Model with Multiple Mediators**

471 Let  $\mathbf{M}$  be a set of  $p$  variables,  $M^{(1)}$ ,  $M^{(2)}$ , to  $M^{(p)}$ , each a potential mediator in the causal pathway between  
472  $A$  and  $Y$ . We assume that the ordering of the potential mediators is arbitrary and that  $Y$  is continuous.  
473 Given a dataset of  $n$  individuals, with  $A_i$ ,  $Y_i$ ,  $\mathbf{M}_i$ , and  $q$  covariates  $\mathbf{C}_i$  measured for each subject  $i$ , we can  
474 evaluate the mediating role of  $\mathbf{M}$  with the models

$$475 \quad E[Y_i|A_i, \mathbf{M}_i, \mathbf{C}_i] = \beta_a A_i + \boldsymbol{\beta}_m^T \mathbf{M}_i + \boldsymbol{\beta}_c^T \mathbf{C}_i \quad (1)$$

476 and

$$477 \quad E[\mathbf{M}_i|A_i, \mathbf{C}_i] = \boldsymbol{\alpha}_a A_i + \boldsymbol{\alpha}_c \mathbf{C}_i. \quad (2)$$

478 We refer to these as the outcome and mediator models. Bolded terms distinguish vectors from  
479 scalars. Under certain assumptions, the parameters of this model can be used to derive causal effects of  
480 interest: Namely, in addition to the baseline assumption of temporality, we assume (1) that there is no  
481 unmeasured confounding in the exposure-outcome association after conditioning on  $\mathbf{C}$ , (2) that there is no  
482 unmeasured confounding in the mediator-outcome associations after adjusting for the exposure and  $\mathbf{C}$ , (3)  
483 that there is no unmeasured confounding of the exposure-mediator associations after conditioning on  $\mathbf{C}$ ,  
484 and (4) that the measured confounders of the mediator-outcome associations are not caused by the  
485 exposure (which would make those confounders mediators themselves). In these circumstances only can  
486  $\beta_a$  be interpreted as the natural direct effect of  $A$  on  $Y$ ,  $\boldsymbol{\alpha}_a^T \boldsymbol{\beta}_m$  the natural indirect effect of  $A$  on  $Y$  through  
487  $\mathbf{M}$ , and  $\beta_a + \boldsymbol{\alpha}_a^T \boldsymbol{\beta}_m$  the total effect of  $A$  on  $Y$ <sup>33</sup>. We say a mediator  $M^{(j)}$  is *active* if  $(\boldsymbol{\alpha}_a)_j (\boldsymbol{\beta}_m)_j$  is not zero,  
488 since it contributes mathematically to the indirect effect, but this contribution itself cannot be formally  
489 interpreted causally unless the mediators are independent conditional on  $A$  and  $\mathbf{C}$ . Extensions of this  
490 framework cover cases when  $Y$  is binary, when  $\mathbf{M}$  is binary, or when the outcome model requires an  
491 interaction effect between  $\mathbf{M}$  and  $A$ <sup>33</sup>.

492 A summary of the methods that can evaluate  $\mathbf{M}$  as a mediator is provided in Table 4, using the  
493 above pair of models as a frame of reference. We describe each of the methods in greater detail in the  
494 following three sections.

495



Name and Author	Estimation of global indirect effect	Estimation of mediation contributions	Mediator identification	Y Data Type	Summary
<b>Group 1 Methods</b>					
HIMA; Zhang, 2016	Point estimation	Point estimation	Yes	Continuous or binary	Fits the outcome model with the minimax concave penalty. Requires subsequent fitting of ordinary least squares regression to test the statistical significance the mediation contributions.
HDMA; Gao, 2019	Point estimation	Point, interval estimation	Yes	Continuous or binary	Fits the outcome model with the de-sparsified LASSO penalty.
MedFix; Zhang, 2021	Point estimation	Point, interval estimation	Yes	Continuous	Fits the outcome model with the adaptive LASSO penalty. Can also be applied when the exposure is high-dimensional in addition to the mediators.
Pathway LASSO Zhao and Luo, 2022	Point estimation	Point estimation	Yes	Continuous	Fits the outcome model and mediator models with a jointly penalized likelihood, directly applying shrinkage to the mediation contributions $(\alpha_a)(\beta_m)$ .
BSLMM; Song, 2020	Bayesian point, interval estimation	Bayesian interval estimation	Yes	Continuous	Bayesian mixed-model in which the mediator-outcome associations $(\beta_m)_j$ and the exposure-mediator associations $(\alpha_a)_j$ are assumed to independently follow sparse normal distributions.
GMM; Song, 2021	Bayesian point, interval estimation	Bayesian interval estimation	Yes	Continuous	Bayesian mixed-model in which the mediator-outcome associations $(\beta_m)_j$ and the exposure-mediator associations $(\alpha_a)_j$ are assumed to jointly follow a sparse multivariate normal distribution.
<b>Group 3 Methods</b>					
PCMA; Huang and Pan, 2016	Point, interval estimation	No	No	Continuous or binary	Applies principal component analysis on the mediator model residuals, transforming the mediators so they are independent. Can be applied when there is $A-M$ interaction in the outcome model.
SPCMA; Zhao, 2019	Point, interval estimation	No	Identifies whether subsets of the mediators are jointly active	Continuous	Similar to PCMA but applies sparse PCA, resulting in transformed mediators that are more interpretable.
HILMA; Zhou, 2020	Point, interval estimation	No	No	Continuous	Uses a debiased penalized regression approach to directly estimate the global indirect effect $\alpha_a^T \beta_m$ . Can be applied for multiple exposures simultaneously.
<b>Group 3 Methods</b>					
HDMM; Chen, 2018	No	No	Nonspecifically identifies groups of active mediators	Continuous	Estimates “directions of mediation” by which the observed mediators can be linearly combined to form latent mediators. The latent mediators replace the true mediators in the analysis.
LVMA; Derkach, 2019	No	No	Identifies inputted mediators associated with latent mediators	Continuous or binary	Reformulates the causal structure of the mediation problem. Assumes that $M$ itself is not responsible for mediation, but rather that the effect of $A$ on $Y$ is mediated by latent, unmeasured factors, $F$ , which also cause changes in $M$ .

## 498 **Group 1 Methods**

499 This group of methods can estimate both the global indirect effect  $\alpha_a^T \beta_m$  and the mediator-specific  
500 contributions  $(\alpha_a)_j(\beta_m)_j, j$  from 1 to  $p$ .

501

### 502 **HIMA**

503 High-dimensional mediation analysis (HIMA), proposed by Zhang et al. (2016), is a penalized regression  
504 approach with two main steps: First, the outcome model is fitted with a minimax concave penalty<sup>54</sup>,  
505 performing feature selection on the mediators by setting some of them to have no effect on  $Y$ <sup>34</sup>. Then,  
506 among the remaining mediators, they fit the mediator models individually using ordinary regression. The  
507 authors test the significance of  $(\alpha_a)_j(\beta_m)_j$  by applying Bonferroni correction to the maximum of the  $(\beta_m)_j$   
508 and  $(\alpha_a)_j$  p-values. To obtain p-values for the  $(\beta_m)_j$  estimates, the authors re-fit the reduced outcome  
509 model by ordinary least squares, which statistically may be overconfident. The authors also recommend  
510 an initial screening step to reduce the number of mediators at the start, as the outcome model will still be  
511 unstable if  $p$  is extremely large compared to  $n$ .

512

### 513 **HDMA**

514 High-dimensional mediation analysis (HDMA), proposed by Gao et al. (2019), is the same as HIMA  
515 except for its penalty function, replacing the minimax concave penalty with the recently-proposed de-  
516 sparsified LASSO<sup>35,55</sup>. The advantage of this penalty is that the resulting estimates of  $\beta_m$  are  
517 asymptotically normal, so one can test their statistical significance without needing to subsequently apply  
518 ordinary least squares. HDMA is also less biased than HIMA when the mediators are highly-correlated.

519

### 520 **MedFix**

521 Mediation analysis via fixed effect model (MedFix) is another extension of HIMA, proposed by Zhang  
522 (2021)<sup>36</sup>. MedFix was originally proposed for a setting where there are not only multiple mediators, but  
523 also multiple exposures, which it handles by applying adaptive LASSO to both the outcome model and  
524 the mediator models. If there is only one exposure, feature selection in the mediator models is not  
525 necessary, and applying MedFix is analogous to applying HDMA except with adaptive LASSO instead of  
526 debiased LASSO.

527

### 528 **Pathway LASSO**

529 Pathway LASSO is another penalized regression approach, proposed by Zhao and Luo (2022)<sup>37</sup>. Whereas  
530 HIMA, HDMA, and MedFix use a two-step design—the outcome model and mediator models fitted  
531 separately—this method fits the models all together, with a jointly penalized likelihood. The penalty not  
532 only applies shrinkage to the mediator-outcome associations, like the other methods, but also to the  
533 exposure-mediator associations and the mediation contributions.

534

### 535 **BSLMM**

536 The Bayesian sparse linear mixed model (BSLMM) is a Bayesian approach proposed by Song et al.  
537 (2020)<sup>15</sup>. The model assumes  $\alpha_a$  and  $\beta_m$  are random vectors, both independently following mixtures of  
538 normal distributions. Most of the effects are presumed to be small, owing to a normal distribution with  
539 mean zero and small variance, while the others are allowed to be larger, resulting from a normal  
540 distribution with higher variance. We estimate the effects with their posterior mean, and we distinguish

541 active mediators from inactive with their posterior inclusion probability of belonging to the distribution  
542 with higher variance.

543

### 544 **GMM**

545 The Gaussian mixed model (GMM), proposed by Song et al. (2021), is an extension of BSLMM in which  
546 the  $(\alpha_a)_j, (\beta_m)_j$  pairs are treated as correlated, following a mixture of multivariate normal distributions  
547 instead of two independent normal distributions<sup>42</sup>. Thus, GMM may be more useful than BSLMM if the  
548 true size of each  $(\beta_m)_j$  is related to the size of the corresponding  $(\alpha_a)_j$ , and vice-versa.

549

## 550 **Group 2 Methods**

551 This group of methods directly estimate the global indirect effect without producing estimates of its  
552 mediator-specific contributions.

553

### 554 **PCMA**

555 Principal component mediation analysis (PCMA), proposed by Huan and Pan (2016), was an early  
556 method for multiple-mediator mediation using principal component analysis (PCA)<sup>39</sup>. The authors  
557 perform PCA on the residual matrix of the mediator models, then use the  $p$  by  $r$  loading matrix  $\mathbf{Q}$  to  
558 transform the matrix  $\mathbf{M}$  into a new set of mediators,  $\mathbf{M}^*$ , which are uncorrelated conditional on  $\mathbf{A}$  and  $\mathbf{C}$ .  
559 The transformed mediators then replace the original mediators in the analysis, and because they are  
560 uncorrelated, the outcome and mediator models can be fit without issue. Although the mediators have  
561 been transformed, and the mediator-specific contributions  $(\alpha_a)_j, (\beta_m)_j$  no longer correspond to the original  
562  $j^{\text{th}}$  mediator, the global indirect effect  $\alpha_a^T \beta_m$  can still be estimated with its original interpretation. The  
563 authors set  $r$  to equal  $p$ , though this is only possible if  $p$  is less than  $n$ .

564

### 565 **SPCMA**

566 Zhao et al (2019) proposed sparse principal component analysis (SPCMA) to improve the interpretability  
567 of the results from PCMA<sup>40</sup>. In PCMA, the transformed mediators are difficult to interpret because they  
568 are sums of all  $p$  original mediators; whereas in SPCMA, the loading matrix  $\mathbf{Q}$  is sparsified, meaning that  
569 each transformed mediator is only a sum of a few of the original mediators. The results are easier to  
570 interpret because, if a specific transformed mediator has a large effect, it can potentially be traced back to  
571 the original mediators which were used to construct it. SPCMA induces bias in its estimation compared to  
572 PCMA, but it can be helpful for identifying groups of mediators which may be active.

573

### 574 **HILMA**

575 High-dimensional linear mediation analysis (HILMA), proposed by Zhou (2020), estimates  $\alpha_a^T \beta_m$  with a  
576 complex, de-biased penalized regression approach<sup>38</sup>. The mathematics of the procedure are beyond the  
577 scope of this text, but the proposed estimator has asymptotic properties for testing whether  $\alpha_a^T \beta_m$  is zero,  
578 and can also be applied when there are multiple (but not high-dimensional) exposures.

579

## 580 **Group 3 Methods**

581 The last group of methods is fundamentally distinct from the others: Instead of fitting the original  
582 mediation models (Group 1), or estimating the mediation effect without fitting the models (Group 2), they  
583 reconceptualize the causal structure of the problem to produce results with unique interpretations. Like

584 any method, they should only be applied when their assumptions about the causal structure are  
585 reasonable.

586

### 587 **HDMM**

588 High-dimensional multivariate mediation (HDMM), proposed by Chén et al. (2018), is similar to PCMA  
589 in that it uses dimension reduction, but chooses the loading vectors with a likelihood-based approach  
590 instead of PCA<sup>41</sup>. The loading vectors are referred to as “directions of mediation,” each vector specifying  
591 a linear combination of mediators which contribute to the likelihood of the mediation models. Hence,  
592 HDMM implicitly assumes that there are latent, unmeasured mediating variables that can be represented  
593 as linear combinations of the observed mediators. The results of HDMM are difficult to interpret, but it  
594 can still be useful for identifying whether there is any mediation through  $M$  at all, and for identifying  
595 large subsets of mediators that contribute to that mediation.

596

### 597 **LVMA**

598 Latent variable mediation analysis (LVMA), proposed by Derkach et al. (2019), assumes that  $M$  itself is  
599 not involved in mediation, but rather, that there are a small number of unmeasured mediators,  $F$ , which  
600 transmit the effect of  $A$  to  $Y$  and which also cause changes in  $M$ <sup>43</sup>. In other words, LVMA assumes  
601 explicitly what HDMM assumes implicitly, and the results of the two methods have a similar structure. A  
602 key feature of LVMA is that the  $F \rightarrow M$  associations are sparsified, meaning that the method can be used  
603 for detecting relevant mediators in  $M$ . An observed mediator would be considered active if it is associated  
604 with a latent mediator that is itself associated with  $A$  and  $Y$ .

605

## 606 **Simulation study**

### 607 **Simulation settings**

608 We evaluate the above methods with a simulation study. To contrast them under diverse  
609 conditions, we consider three different settings of mediation: (1) a baseline setting in which the mediation  
610 signals are sparse and the (potential) mediators are moderately correlated, (2) a high-correlation setting  
611 with sparse signals, and (3) a moderate correlation setting in which the signals are non-sparse. Within  
612 each of these settings, we also vary the degree of mediation by modifying three parameters: the  
613 proportion of variance in  $M$  that is explained by  $A$  among those associated with  $A$  ( $PVE_A$ ), the proportion  
614 of the variance of  $Y$  that is explained by the direct effect ( $PVE_{DE}$ ), and the proportion of the variance of  $Y$   
615 that is explained by the global indirect effect ( $PVE_{IE}$ ). For a baseline case, we let  $PVE_A$  equal 0.20 and  
616  $PVE_{DE}$  and  $PVE_{IE}$  both equal 0.10; then, in three additional cases, we sequentially decrease one of these  
617 parameters by half, weakening the signal, and set the other two parameters to their values from the  
618 baseline. Between Settings (1) to (3), this amounted to 12 unique data-generating mechanisms in total.  
619 Each of these was evaluated with a sample size of 1,000 and 2,500, with the number of potential  
620 mediators fixed at 2,000. All combinations of settings are listed below in Table 5.

621

622

623

624

625

626

627

628  
629  
630

**Table 5. Complete list of settings in simulation study**

Number of potential mediators ( $p$ )	Sample Size ( $n$ )	Sparsity of signals	Degree of correlation	PVE <sub>A</sub>	PVE <sub>IE</sub>	PVE <sub>DE</sub>
2000	2500	Sparse	Baseline	0.20	0.10	0.10
2000	2500	Sparse	Baseline	0.20	0.05	0.10
2000	2500	Sparse	Baseline	0.10	0.10	0.10
2000	2500	Sparse	Baseline	0.20	0.10	0.05
2000	2500	Sparse	High	0.20	0.10	0.10
2000	2500	Sparse	High	0.20	0.05	0.10
2000	2500	Sparse	High	0.10	0.10	0.10
2000	2500	Sparse	High	0.20	0.10	0.05
2000	2500	Non-sparse	Baseline	0.20	0.10	0.10
2000	2500	Non-sparse	Baseline	0.20	0.05	0.10
2000	2500	Non-sparse	Baseline	0.10	0.10	0.10
2000	2500	Non-sparse	Baseline	0.20	0.10	0.05
2000	1000	Sparse	Baseline	0.20	0.10	0.10
2000	1000	Sparse	Baseline	0.20	0.05	0.10
2000	1000	Sparse	Baseline	0.10	0.10	0.10
2000	1000	Sparse	Baseline	0.20	0.10	0.05
2000	1000	Sparse	High	0.20	0.10	0.10
2000	1000	Sparse	High	0.20	0.05	0.10
2000	1000	Sparse	High	0.10	0.10	0.10
2000	1000	Sparse	High	0.20	0.10	0.05
2000	1000	Non-sparse	Baseline	0.20	0.10	0.10
2000	1000	Non-sparse	Baseline	0.20	0.05	0.10
2000	1000	Non-sparse	Baseline	0.10	0.10	0.10
2000	1000	Non-sparse	Baseline	0.20	0.10	0.05

631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649

**Simulated dataset creation**

First, to obtain sparse mediation effects for Settings (1) and (2), we assume that 1,920 of the 2,000 coefficients  $(\alpha_a)_j$  and  $(\beta_m)_j$  are zero and the remaining 80 are standard normal. Twenty of the nonzero  $(\alpha_a)_j$  and  $(\beta_m)_j$  are chosen to overlap and have  $(\alpha_a)_j(\beta_m)_j$  not equal to zero. To obtain non-sparse signals for Setting (3), we sample the previously zero coefficients from a normal distribution with mean zero and standard deviation 0.2. (These parameter vectors are sampled only once, at the start of the simulations, so that the global mediation effect is held constant, but we shuffle the mediators in each dataset so that different mediators are assigned the effects each time.) Once we have these, we obtain a single simulated dataset by sampling  $A_i$  from a standard normal distribution, then produce  $M_i$  from model (4) assuming there are no covariates. We add noise to  $M_i$  by sampling residuals from a multivariate normal distribution with mean  $\mathbf{0}_p$  and variance  $\Sigma$ , where  $\Sigma$  is derived by shuffling and then tuning the variance-covariance of the observed methylation data (see supplementary section 1). In Settings (1) and (3), we tune  $\Sigma$  so that the correlations between mediators range from -0.37 to 0.49, and in Setting (2), so that they range from -0.58 to 0.75. We fix PVE<sub>A</sub> by scaling  $\Sigma$  appropriately based on  $\alpha_a$ . Finally, we define  $Y_i$  based on model (3) assuming the residuals are Normal(0,  $\sigma^2$ ), choosing  $\beta_a$  and  $\sigma^2$  to yield the desired PVE<sub>DE</sub> and PVE<sub>IE</sub>.

## 650 Evaluation

651 We evaluate the methods by applying them to 100 replicates of each setting in Table 5. We omit  
652 SPCMA, GMM, and LVMA for computational reasons, as they are too computationally costly to deploy  
653 on so many replicates, and omit HDMM because it does not have an estimand that is comparable to the  
654 others. We include a one-at-a-time approach—in which the mediator are assessed individually using  
655 traditional mediation analysis and the joint significance test<sup>44</sup>—as a baseline for comparison. When  
656 running HIMA, HDMA, MedFix, and pathway LASSO, we pre-screen the mediators to only include the  
657  $n/\log(n)$  mediators with the strongest associations with  $Y$  adjusting for  $A$ , which is the approach  
658 recommended by the HIMA and HDMA authors<sup>34,35</sup> (see supplementary section 2 for more details). For  
659 comparison metrics, we use (1) the true positive rate for detecting active mediators,  $\text{TPR} =$   
660  $\frac{\text{number of true mediators identified}}{\text{number of true mediators}}$ ; (2) the mean squared error in estimating the mediation contributions of  
661 inactive mediators,  $\text{MSE}_{\text{Inactive}} = \text{mean}_{j: \text{Inactive}} \left( (\widehat{\alpha}_a)_j (\widehat{\beta}_m)_j - (\alpha_a)_j (\beta_m)_j \right)^2$ ; (3) the mean squared  
662 error in estimating the mediation contributions of active mediators,  $\text{MSE}_{\text{Active}} =$   
663  $\text{mean}_{j: \text{Active}} \left( (\widehat{\alpha}_a)_j (\widehat{\beta}_m)_j - (\alpha_a)_j (\beta_m)_j \right)^2$ ; and (4) the percent relative bias in estimating the global  
664 indirect effect,  $\frac{|\widehat{\alpha}_a^T \widehat{\beta}_m - \alpha_a^T \beta_m|}{\alpha_a^T \beta_m} \times 100$ . In the non-sparse setting, since all the mediators contribute to the  
665 indirect effect, we consider the “active” ones to be those whose mediator-outcome and exposure-mediator  
666 effects both come from the distribution with higher-variance, and the others inactive. Each metric is  
667 computed for each dataset to the applicable methods, and we report the average and a 95% empirical  
668 confidence interval over the 100 replicates.

## 669 Scalability comparison

671 We compare the scalability of the methods by assessing their processing time on simulated  
672 datasets of two sizes: one with 100 observations and 200 mediators and one with 1,000 observations and  
673 2,000 mediators. For the larger dataset, we use one of the datasets created for the simulation study, and  
674 for the smaller dataset, we subset the rows and columns of  $M$  and the entries in  $A$  and  $Y$ . Run times are  
675 assessed on a single core of an Intel(R) Xeon(R) Gold 6242R CPU @ 3.10GHz processor. We attempt  
676 each method 30 times and report the mean and interquartile range of the computation times. Since  
677 SCPMA and BSLMM tend to be time-consuming, we approximate their run times by downscaling the  
678 appropriate parameters: In particular, since the desired number of principal components in SPCMA is  
679 100, we use only 2 principal components and scale the computing time by 50; and since the desired  
680 number of posterior samples in BSLMM is 30,000, we draw only 750 samples and scale the result by 40.  
681 Ad hoc experimentation confirmed that the methods were approximately linear with respect to these  
682 inputs.

683

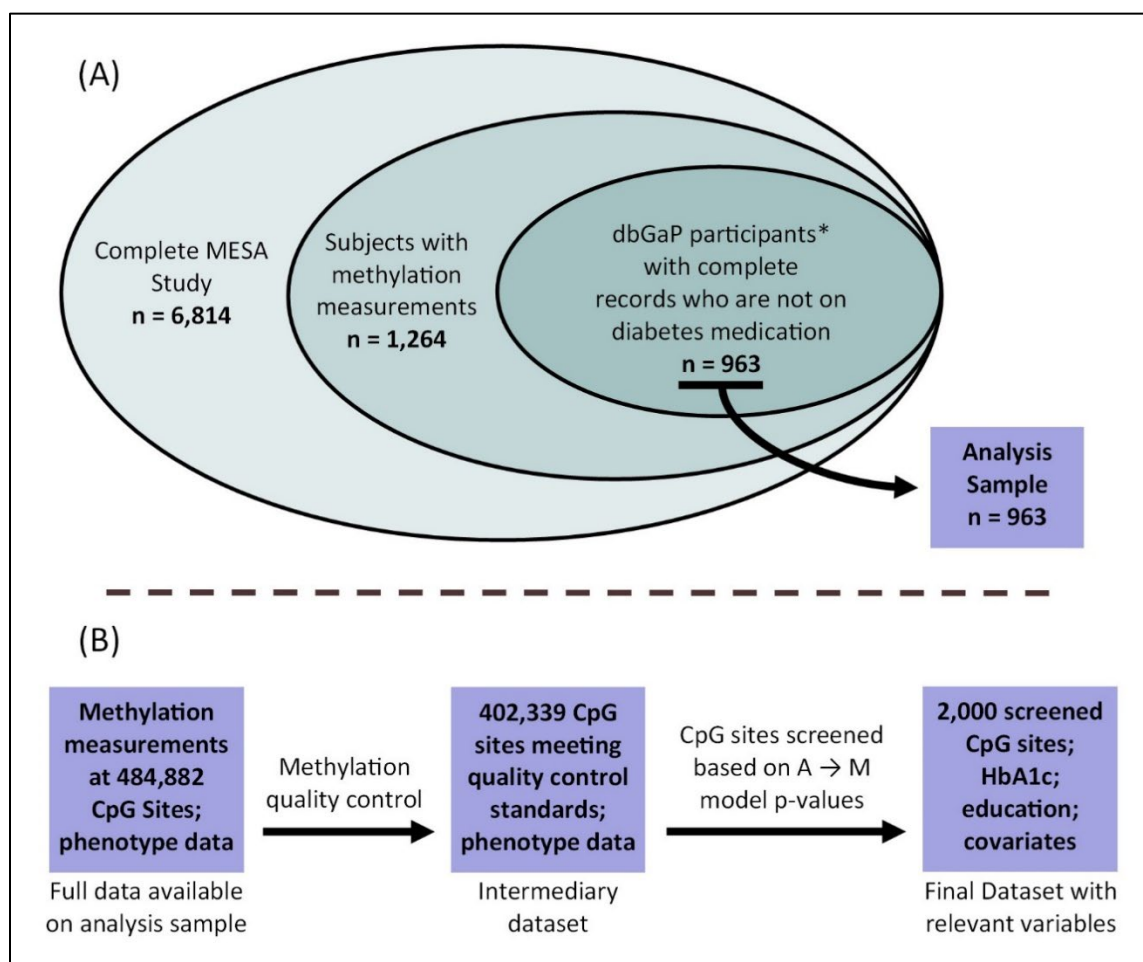
## 684 Data application with MESA

685 To demonstrate how these methods can be applied to observed DNAm data, we evaluate the  
686 association between SES and HbA1c and its potential mediation through DNAm. For the exposure, we  
687 use a binary variable that indicates low educational attainment (less than a 4-year college degree); for the  
688 outcome, we use HbA1c, a continuous variable that reflects average three-month blood glucose level. Our  
689 data for this portion come from the Multi-Ethnic Study of Atherosclerosis (MESA), a US population-  
690 based longitudinal study<sup>17</sup>. Out of 6,814 total participants, a random subsample of 1,264 had their DNAm

691 measured at 484,882 CpG sites. We limit our analysis to the 963 participants who (1) had methylation  
 692 data, (2) had no missing data for the required variables, (3) consented to genetic and phenotypic use  
 693 through the database of Genotypes and Phenotypes (dbGaP) (phs000209.v13.p3), and (4) were not on  
 694 diabetes medication, which can cause changes in HbA1c (Fig. 6). Standard quality control filters reduced  
 695 the number of CpG sites to 402,339. Since it is not statistically or computationally feasible to include so  
 696 many mediators at once, we used a screening procedure to reduce that number further, fitting model (6)  
 697 below for each mediator separately to choose the 2,000 CpG sites at which DNAm was most strongly  
 698 associated with education based on the  $(\alpha_a)_j$  p-value. These 2,000 formed the baseline set of CpGs for our  
 699 analysis. DNAm was measured using M-values, defined as the log-2 ratio of the methylated to  
 700 unmethylated probe intensities, which has the advantage of occurring on a continuous and unbounded  
 701 scale<sup>56</sup>. For more details see supplementary section 3. A model for the proposed mechanism is given by  
 702  $E[\text{HbA1c}_i | \text{Education}_i, \text{DNAm}_i, \text{Covariates}_i] = \beta_a \text{Education}_i + \beta_m^T \text{DNAm}_i + \beta_c^T \text{Covariates}_i$  (5)  
 703 and

$$E[\text{DNAm}_i | \text{Education}_i, \text{Covariates}_i] = \alpha_a \text{Education}_i + \alpha_c \text{Covariates}_i, \quad (6)$$

704 where the covariates include age, sex, race, and the estimated proportions of residual non-monocytes (i.e.,  
 705 neutrophils, B cells, T cells, and natural killer cells) as fixed effects and methylation chip and position as  
 706 random effects.  
 707



708 **Fig. 6. Pre-processing of MESA methylation data.** \*Participants who consent to genetic and phenotypic data use,  
 709 and whose data is available on dbGaP.

710 We performed mediation analysis on the final dataset of 963 individuals and 2,000 CpG sites. All of the  
711 mediation methods described above were included except for GMM and LVMA, which again were too  
712 costly computationally. Although it is reasonable for some of the methods to include all 2,000 CpG sites  
713 directly in the multivariable model, HIMA and HDMA involve sure independence screening<sup>57</sup> to reduce  
714 the number of mediators in advance to  $n/\log(n)$ , where  $n$  is the sample size. For the sake of consistency  
715 across the penalized regression methods, we do so with not only HIMA and HDMA, but also MedFix and  
716 pathway LASSO, including only the 141 ( $963/\log(963)$ ) CpG sites most associated with low education (a  
717 direct extension of the initial screening). (Note that, for HIMA and HDMA, this screening is *part* of the  
718 proposed method, not separate from it, but for MedFix and pathway LASSO the additional screening is  
719 still beneficial for the sake of comparing methods and for statistical and computational efficiency).  
720 Additional pre-screening is not necessary for PCMA, SPCMA, BSLMM, and HILMA, and we include all  
721 2,000 CpG sites directly; however, in HDMM, which cannot accommodate  $p > n$  simplistically, we again  
722 use only twice-screened subset of 141 sites. For the sake of comparison with multivariate methods, we  
723 also include a one-at-a-time mediation method based on linear regression and the joint significance test.  
724 For estimating the total effect, the methods PCMA, SPCMA, BSLMM, and Pathway LASSO all produce  
725 estimates of the direct effect, so we can estimate the total effect by summing the estimated direct and  
726 global indirect effects. Since the methods HIMA, HDMA, and MedFix do not produce estimates of the  
727 direct effect, we first estimate the total effect on its own by fitting model (5) with the mediators excluded,  
728 then subtract the estimated global indirect effect from this value to obtain an estimate of the direct effect.  
729 As none of the high-dimensional methods are built to directly handle random effects as covariates, we  
730 regress these out of the outcome variable and potential mediators in advance. For the fixed effect  
731 covariates, HIMA, HDMA, MedFix, and BSLMM allow one to include them directly; whereas in PCMA,  
732 SPCMA, HILMA, HDMM, and pathway LASSO, we regressed them out in advance from the outcome  
733 and mediators. Continuous variables (including HbA1c and the mediators) were standardized for all  
734 methods. All analysis was conducted using R version 4.2.1.

735

## 736 Data Availability

737 Data used for the simulation study are available from the authors upon request. Data used in the DNAm  
738 analysis can be obtained through the MESA Data Coordinating Center (<https://www.mesahlbi.org/>).

739

## 740 Code Availability

741 R scripts for the analysis are available at [https://github.com/dclarkboucher/mediation\\_DNAm](https://github.com/dclarkboucher/mediation_DNAm). Our R  
742 package “hdmed” can be found at <https://github.com/dclarkboucher/hdmed>.

743

## 744 References

- 745 1. Moore, L. D., Le, T. & Fan, G. DNA methylation and its basic function. *Neuropsychopharmacol.*  
746 *Off. Publ. Am. Coll. Neuropsychopharmacol.* **38**, 23–38 (2013).
- 747 2. Kurdyukov, S. & Bullock, M. DNA Methylation Analysis: Choosing the Right Method. *Biology*  
748 *(Basel)*. **5**, (2016).
- 749 3. Dick, K. J. *et al.* DNA methylation and body-mass index: a genome-wide analysis. *Lancet*  
750 *(London, England)* **383**, 1990–1998 (2014).
- 751 4. Wahl, S. *et al.* Epigenome-wide association study of body mass index, and the adverse outcomes  
752 of adiposity. *Nature* **541**, 81–86 (2017).
- 753 5. Volkmar, M. *et al.* DNA methylation profiling identifies epigenetic dysregulation in pancreatic



- 754 islets from type 2 diabetic patients. *EMBO J.* **31**, 1405–1426 (2012).
- 755 6. Chilunga, F. P. *et al.* Genome-wide DNA methylation analysis on C-reactive protein among  
756 Ghanaians suggests molecular links to the emerging risk of cardiovascular diseases. *NPJ genomic*  
757 *Med.* **6**, 46 (2021).
- 758 7. Nakatochi, M. *et al.* Epigenome-wide association of myocardial infarction with DNA methylation  
759 sites at loci related to cardiovascular disease. *Clin. Epigenetics* **9**, 54 (2017).
- 760 8. Fujii, R. *et al.* Dietary fish and  $\omega$ -3 polyunsaturated fatty acids are associated with leukocyte  
761 ABCA1 DNA methylation levels. *Nutrition* **81**, 110951 (2021).
- 762 9. Sun, Y. V *et al.* Epigenomic association analysis identifies smoking-related DNA methylation  
763 sites in African Americans. *Hum. Genet.* **132**, 1027–1037 (2013).
- 764 10. Philibert, R. A., Plume, J. M., Gibbons, F. X., Brody, G. H. & Beach, S. R. H. The impact of  
765 recent alcohol use on genome wide DNA methylation signatures. *Front. Genet.* **3**, 54 (2012).
- 766 11. Rider, C. F. & Carlsten, C. Air pollution and DNA methylation: effects of exposure in humans.  
767 *Clin. Epigenetics* **11**, 131 (2019).
- 768 12. Lam, L. L. *et al.* Factors underlying variable DNA methylation in a human community cohort.  
769 *Proc. Natl. Acad. Sci. U. S. A.* **109 Suppl**, 17253–17260 (2012).
- 770 13. Needham, B. L. *et al.* Life course socioeconomic status and DNA methylation in genes related to  
771 stress reactivity and inflammation: The multi-ethnic study of atherosclerosis. *Epigenetics* **10**,  
772 958–969 (2015).
- 773 14. Fujii, R., Sato, S., Tsuboi, Y., Cardenas, A. & Suzuki, K. DNA methylation as a mediator of  
774 associations between the environment and chronic diseases: A scoping review on application of  
775 mediation analysis. *Epigenetics* 1–27 (2021) doi:10.1080/15592294.2021.1959736.
- 776 15. Song, Y. *et al.* Bayesian shrinkage estimation of high dimensional causal mediation effects in  
777 omics studies. *Biometrics* **76**, 700–710 (2020).
- 778 16. Du, J. *et al.* Methods for Large-scale Single Mediator Hypothesis Testing: Possible Choices and  
779 Comparisons. (2022) doi:10.48550/arxiv.2203.13293.
- 780 17. Bild, D. E. *et al.* Multi-Ethnic Study of Atherosclerosis: objectives and design. *Am. J. Epidemiol.*  
781 **156**, 871–881 (2002).
- 782 18. Whitaker, S. M. *et al.* The Association Between Educational Attainment and Diabetes Among  
783 Men in the United States. *American journal of men's health* vol. 8 (2014).
- 784 19. Sakurai, M. *et al.* HbA1c and the risks for all-cause and cardiovascular mortality in the general  
785 Japanese population: NIPPON DATA90. *Diabetes Care* **36**, 3759–3765 (2013).
- 786 20. Singer, D. E., Nathan, D. M., Anderson, K. M., Wilson, P. W. & Evans, J. C. Association of  
787 HbA1c with prevalent cardiovascular disease in the original cohort of the Framingham Heart  
788 Study. *Diabetes* **41**, 202–208 (1992).
- 789 21. Yeung, S. L. A., Luo, S. & Schooling, C. M. The Impact of Glycated Hemoglobin (HbA(1c)) on  
790 Cardiovascular Disease Risk: A Mendelian Randomization Study Using UK Biobank. *Diabetes*  
791 *Care* **41**, 1991–1997 (2018).
- 792 22. Borghol, N. *et al.* Associations with early-life socio-economic position in adult DNA methylation.  
793 *Int. J. Epidemiol.* **41**, 62–74 (2012).
- 794 23. Chen, Z. *et al.* DNA methylation mediates development of HbA1c-associated complications in  
795 type 1 diabetes. *Nat. Metab.* **2**, 744–762 (2020).
- 796 24. Baron, R. M. & Kenny, D. A. The Moderator-Mediator Variable Distinction in Social  
797 Psychological Research: Conceptual, Strategic, and Statistical Considerations. *Journal of*  
798 *personality and social psychology* vol. 51.
- 799 25. MacKinnon, D. *Introduction to statistical mediation analysis.* (New York, NY u.a: Erlbaum).
- 800 26. VanderWeele, T. J. Marginal Structural Models for the Estimation of Direct and Indirect Effects.  
801 *Epidemiology* **20**, (2009).
- 802 27. Pearl, J. Direct and Indirect Effects. in *Proceedings of the Seventeenth Conference on Uncertainty*  
803 *in Artificial Intelligence* 411–420 (Morgan Kaufmann Publishers Inc., 2001).
- 804 28. Robins, J. M. & Greenland, S. Identifiability and exchangeability for direct and indirect effects.

- 805 *Epidemiology* **3**, 143–155 (1992).
- 806 29. VanderWeele, T. J. Mediation Analysis: A Practitioner’s Guide. *Annu. Rev. Public Health* **37**, 17–  
807 32 (2016).
- 808 30. VanderWeele author., T. *Explanation in causal inference: methods for mediation and interaction*.  
809 *Explanation in causal inference: methods for mediation and interaction* (Oxford University Press,  
810 2015).
- 811 31. Du, J. *et al.* Methods for large-scale single mediator hypothesis testing: Possible choices and  
812 comparisons. *Genet. Epidemiol.* **n/a**, (2022).
- 813 32. Aung, M. T. *et al.* Application of an analytical framework for multivariate mediation analysis of  
814 environmental data. *Nat. Commun.* **11**, 5624 (2020).
- 815 33. VanderWeele, T. J. & Vansteelandt, S. Mediation Analysis with Multiple Mediators. *Epidemiol.*  
816 *Method.* **2**, 95–115 (2014).
- 817 34. Zhang, H. *et al.* Estimating and testing high-dimensional mediation effects in epigenetic studies.  
818 *Bioinformatics* **32**, 3150–3154 (2016).
- 819 35. Gao, Y. *et al.* Testing Mediation Effects in High-Dimensional Epigenetic Studies. *Front. Genet.*  
820 **10**, 1195 (2019).
- 821 36. Zhang, Q. High-Dimensional Mediation Analysis with Applications to Causal Gene Identification.  
822 *Statistics in biosciences* (2021).
- 823 37. Zhao, Y. & Luo, X. Pathway LASSO: pathway estimation and selection with high-dimensional  
824 mediators. *Stat. Interface* **15**, 39–50 (2022).
- 825 38. Zhou, R. R., Wang, L. & Zhao, S. D. Estimation and inference for the indirect effect in high-  
826 dimensional linear mediation models. *Biometrika* **107**, 573–589 (2020).
- 827 39. Huang, Y.-T. & Pan, W.-C. Hypothesis test of mediation effect in causal mediation model with  
828 high-dimensional continuous mediators. *Biometrics* **72**, 402–413 (2016).
- 829 40. Zhao, Y., Lindquist, M. A. & Caffo, B. S. Sparse principal component based high-dimensional  
830 mediation analysis. *Comput. Stat. Data Anal.* **142**, 106835 (2020).
- 831 41. Chén, O. Y. *et al.* High-dimensional multivariate mediation with application to neuroimaging data.  
832 *Biostatistics (Oxford, England)* vol. 19 (2018).
- 833 42. Song, Y. *et al.* Bayesian sparse mediation analysis with targeted penalization of natural indirect  
834 effects. *J. R. Stat. Soc. Ser. C (Applied Stat.)* **70**, 1391–1412 (2021).
- 835 43. Derkach, A., Pfeiffer, R. M., Chen, T.-H. & Sampson, J. N. High dimensional mediation analysis  
836 with latent variables. *Biometrics* **75**, 745–756 (2019).
- 837 44. MacKinnon, D. P., Lockwood, C. M., Hoffman, J. M., West, S. G. & Sheets, V. A comparison of  
838 methods to test mediation and other intervening variable effects. *Psychol. Methods* **7**, 83–104  
839 (2002).
- 840 45. Pedroso, J. A. B., Ramos-Lobo, A. M. & Donato, J. J. SOCS3 as a future target to treat metabolic  
841 disorders. *Hormones (Athens)*. **18**, 127–136 (2019).
- 842 46. Wang, Y. Z. *et al.* DNA Methylation Mediates the Association Between Individual and  
843 Neighborhood Social Disadvantage and Cardiovascular Risk Factors. *Front. Cardiovasc. Med.* **9**,  
844 848768 (2022).
- 845 47. Stage, E. *et al.* The effect of the top 20 Alzheimer disease risk genes on gray-matter density and  
846 FDG PET brain metabolism. *Alzheimer’s Dement. (Amsterdam, Netherlands)* **5**, 53–66 (2016).
- 847 48. Mei, H. *et al.* Tissue Non-Specific Genes and Pathways Associated with Diabetes: An Expression  
848 Meta-Analysis. *Genes (Basel)*. **8**, (2017).
- 849 49. Rahman, S. A., Nessa, A. & Hussain, K. Molecular mechanisms of congenital hyperinsulinism. *J.*  
850 *Mol. Endocrinol.* **54**, R119–R129 (2015).
- 851 50. Galcheva, S., Al-Khawaga, S. & Hussain, K. Diagnosis and management of hyperinsulinaemic  
852 hypoglycaemia. *Best Pract. Res. Clin. Endocrinol. Metab.* **32**, 551–573 (2018).
- 853 51. Pedroso, J. A. B. *et al.* Inactivation of SOCS3 in leptin receptor-expressing cells protects mice  
854 from diet-induced insulin resistance but does not prevent obesity. *Mol. Metab.* **3**, 608–618 (2014).
- 855 52. Senniappan, S., Shanti, B., James, C. & Hussain, K. Hyperinsulinaemic hypoglycaemia: genetic

- 856 mechanisms, diagnosis and management. *J. Inherit. Metab. Dis.* **35**, 589–601 (2012).  
857 53. Huang, Y.-T., Vanderweele, T. J. & Lin, X. Joint analysis of SNP and gene expression data in  
858 genetic association studies of complex diseases. *Ann. Appl. Stat.* **8**, 352–376 (2014).  
859 54. Zhang, C.-H. Nearly unbiased variable selection under minimax concave penalty. *Ann. Stat.* **38**,  
860 894–942 (2010).  
861 55. Zhang, S. S. & Zhang, C.-H. Confidence intervals for low dimensional parameters in high  
862 dimensional linear models. *Journal of the Royal Statistical Society. Series B, Statistical*  
863 *methodology* vol. 76 (2014).  
864 56. Du, P. *et al.* Comparison of Beta-value and M-value methods for quantifying methylation levels  
865 by microarray analysis. *BMC Bioinformatics* **11**, 587 (2010).  
866 57. Fan, J. & Lv, J. Sure independence screening for ultrahigh dimensional feature space. *J. R. Stat.*  
867 *Soc.* **70**, 849–911 (2008).  
868

## 869 **Acknowledgements**

870 MESA and the MESA SHARe project are conducted and supported by the National Heart, Lung, and  
871 Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA is provided by  
872 contracts 75N92020D00001, HHSN268201500003I, N01-HC-95159, 75N92020D00005, N01-HC-  
873 95160, 75N92020D00002, N01-HC-95161, 75N92020D00003, N01-HC-95162, 75N92020D00006, N01-  
874 HC 95163, 75N92020D00004, N01-HC-95164, 75N92020D00007, N01-HC-95165, N01-HC- 95166,  
875 N01-HC-95167, N01-HC-95168, N01-HC-95169, UL1-TR-000040, UL1-TR-001079, UL1-TR-001420,  
876 UL1-TR-001881, and DK063491. The MESA Epigenomics & Transcriptomics Studies were funded by  
877 NIH grants 1R01HL101250, 1RF1AG054474, R01HL126477, R01DK101921,  
878 and R01HL135009. Co-authors of this manuscripts were partially supported by NHLBI grant  
879 R01HL141292, NSF grant DMS1712933, and NIH grants R01HG008773 and 1UG3CA267907.