

Web-based Supporting Materials for Bayesian Sparse Mediation Analysis with Targeted Penalization of Natural Indirect Effects

Yanyi Song¹, Xiang Zhou^{1,*}, Jian Kang^{1,**}, Max T. Aung¹, Min Zhang¹, Wei Zhao²,
Belinda L. Needham², Sharon L. R. Kardia², Yongmei Liu³, John D. Meeker⁴, Jennifer A. Smith²,
and Bhramar Mukherjee¹

¹Department of Biostatistics, University of Michigan, Ann Arbor, MI, U.S.A.

²Department of Epidemiology, University of Michigan, Ann Arbor, MI, U.S.A.

³Division of Cardiology, Department of Medicine, Duke University School of Medicine, Durham, NC, U.S.A.

⁴Department of Environmental Health Sciences, University of Michigan, Ann Arbor, MI, U.S.A.

* *email*: xzhousph@umich.edu

** *email*: jiankang@umich.edu

1 Identifiability Assumptions for Causal Mediation Analysis

We use the same counterfactual notation as in the main manuscript. To connect potential variables to observed data, we make the Stable Unit Treatment Value Assumption (SUTVA) [1, 2]. Specifically, the SUTVA assumes there is no interference between subjects and the consistency assumption, which states that the observed variables are the same as the potential variables corresponding to the actually observed treatment level, i.e., $\mathbf{M}_i = \sum_a \mathbf{M}_i(a)I(A_i = a)$, and $Y_i = \sum_a \sum_m Y_i(a, \mathbf{m})I(A_i = a, \mathbf{M}_i = \mathbf{m})$, where $I(\cdot)$ is the indicator function.

Causal effects are formally defined in terms of potential variables which are not necessarily observed, but the identification of causal effects must be based on observed data. Therefore further assumptions regarding the confounders are required for the identification of causal effects in mediation analysis [3]. We will use $A \perp\!\!\!\perp B | C$ to denote that A is independent of B conditional on C . To estimate the average NDE and NIE from observed data, the following assumptions are needed: (1) $Y_i(a, \mathbf{m}) \perp\!\!\!\perp A_i | \mathbf{C}_i$, no unmeasured confounding for exposure-outcome relationship; (2) $Y_i(a, \mathbf{m}) \perp\!\!\!\perp \mathbf{M}_i | \{\mathbf{C}_i, A_i\}$, no unmeasured confounding for any of mediator-outcome relationship after controlling for the exposure; (3) $\mathbf{M}_i(a) \perp\!\!\!\perp A_i | \mathbf{C}_i$, no un-

measured confounding for the exposure effect on all the mediators; (4) $Y_i(a, \mathbf{m}) \perp\!\!\!\perp \mathbf{M}_i(a^*) | \mathbf{C}_i$, no downstream effect of the exposure that confounds any mediator-outcome relationship. The four assumptions are required to hold with respect to the whole set of mediators. Finally, as in all mediation analysis, the temporal ordering assumption also needs to be satisfied, i.e., the exposure precedes the mediators, and the mediators precede the outcome.

2 Posterior Sampling Algorithm Details for Gaussian Mixture Model (GMM)

Let $\Theta_{GMM} = (\boldsymbol{\beta}_m, \boldsymbol{\alpha}_a, \mathbf{V}_k, \beta_a, \boldsymbol{\beta}_c, \boldsymbol{\alpha}_c, \{\gamma_j\}_{j=1}^p, \pi_k, k = 1, 2, 3, 4, \sigma_e^2, \boldsymbol{\Sigma}, \sigma_a^2)$ denote all the unknown parameters in our Gaussian mixture model. The joint likelihood of $\{Y_i, \mathbf{M}_i\}_{i=1}^n$ given Θ_{GMM} is,

$$\begin{aligned}
\log P(\{Y_i, \mathbf{M}_i\}_{i=1}^n | \Theta_{GMM}, \{A_i, \mathbf{C}_i\}_{i=1}^n) &= \sum_{i=1}^n \log P(Y_i, \mathbf{M}_i | \Theta_{GMM}, A_i, \mathbf{C}_i) \\
&= \sum_{i=1}^n \log P(Y_i | \mathbf{M}_i, \boldsymbol{\beta}_m, \sigma_e^2, \beta_a, \boldsymbol{\beta}_c, A_i, \mathbf{C}_i) \\
&\quad + \log P(\mathbf{M}_i | \boldsymbol{\alpha}_a, \boldsymbol{\alpha}_c, \boldsymbol{\Sigma}, A_i, \mathbf{C}_i) \\
&= \sum_{i=1}^n -\frac{1}{2} \log \sigma_e^2 - \frac{1}{2\sigma_e^2} (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i \beta_a - \mathbf{C}_i^T \boldsymbol{\beta}_c)^2 \\
&\quad - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i)^T \boldsymbol{\Sigma}^{-1} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i)
\end{aligned}$$

The joint log posterior distribution is,

$$\begin{aligned}
\log P(\Theta_{GMM} | \{Y_i, \mathbf{M}_i, A_i, \mathbf{C}_i\}_{i=1}^n) &\propto \sum_{i=1}^n \log P(Y_i | \mathbf{M}_i, \boldsymbol{\beta}_m, \sigma_e^2, \beta_a, \boldsymbol{\beta}_c, A_i, \mathbf{C}_i) + \log P(\mathbf{M}_i | \boldsymbol{\alpha}_a, \boldsymbol{\alpha}_c, \boldsymbol{\Sigma}, A_i, \mathbf{C}_i) \\
&\quad + \log P(\Theta_{GMM}) \\
&= \sum_{i=1}^n -\frac{1}{2} \log \sigma_e^2 - \frac{1}{2\sigma_e^2} (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i \beta_a - \mathbf{C}_i^T \boldsymbol{\beta}_c)^2 \\
&\quad - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i)^T \boldsymbol{\Sigma}^{-1} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i) \\
&\quad + \sum_{j=1}^p \sum_{k=1}^4 \gamma_{jk} \left(-\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{V}_k| - \frac{1}{2} \begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix}^T \mathbf{V}_k^{-1} \begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix} \right) \\
&\quad - \frac{q}{2} \log 2\pi \sigma_c^2 - \frac{\boldsymbol{\beta}_c^T \boldsymbol{\beta}_c}{2\sigma_c^2} - \frac{pq}{2} \log 2\pi \sigma_c^2 - \sum_{j=1}^p \frac{\boldsymbol{\alpha}_{cj}^T \boldsymbol{\alpha}_{cj}}{2\sigma_c^2} \\
&\quad + \sum_{j=1}^p \sum_{k=1}^4 \gamma_{jk} \log(\pi_k) \\
&\quad + \sum_{k=1}^4 a_k \log(\pi_k) + \sum_{k=1}^4 \left(-\frac{\nu + d + 1}{2} \log |\mathbf{V}_k| + \frac{1}{2} \text{tr}(\boldsymbol{\Psi}_0 \mathbf{V}_k^{-1}) \right)
\end{aligned}$$

Sampling $\begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix}$ and γ_{jk}

$$\log p\left(\begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix} \mid \gamma_{jk} = 1, \cdot\right) \propto -\frac{1}{2} \begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix}^T (\mathbf{W}_j + \mathbf{V}_k^{-1}) \begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix} + \mathbf{w}_j^T \begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix}$$

where $\mathbf{W}_j = \begin{bmatrix} \sum_{i=1}^n (\sigma_e^2)^{-1} M_{ij}^2 & 0 \\ 0 & \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} A_i^2 \end{bmatrix}$ ($\boldsymbol{\Sigma}$ is diagonal, and can be replaced as σ_g^2),

and $\mathbf{w}_j = (\sum_{i=1}^n (\sigma_e^2)^{-1} (Y_i - A_i \beta_a - \sum_{j' \neq j} M_{ij'} (\boldsymbol{\beta}_m)_{j'}) M_{ij}, \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} M_{ij} A_i)^T$

$$p\left(\begin{bmatrix} (\boldsymbol{\beta}_m)_j \\ (\boldsymbol{\alpha}_a)_j \end{bmatrix} \mid \gamma_{jk} = 1, \cdot\right) \sim MVN_2((\mathbf{W}_j + \mathbf{V}_k^{-1})^{-1} \mathbf{w}_j, (\mathbf{W}_j + \mathbf{V}_k^{-1})^{-1})$$

$$\log p(\gamma_{jk} = 1 | \cdot) \propto -\frac{1}{2} \log |\mathbf{W}_j \mathbf{V}_k + \mathbf{I}_2| + \frac{1}{2} \mathbf{w}_j^T (\mathbf{W}_j + \mathbf{V}_k^{-1})^{-1} \mathbf{w}_j + \log(\pi_k)$$

Sampling π_k

$$\{\pi_1, \pi_2, \pi_3, \pi_4\} \propto \text{Dirichlet}(a_1 + \sum_{j=1}^p \gamma_{j1}, a_2 + \sum_{j=1}^p \gamma_{j2}, a_3 + \sum_{j=1}^p \gamma_{j3}, a_4 + \sum_{j=1}^p \gamma_{j4})$$

Sampling \mathbf{V}_k

$$\log p(\mathbf{V}_k | \cdot) \propto -\frac{1}{2} \left(\sum_{j=1}^p \gamma_{jk} + \nu + d + 1 \right) \log |\mathbf{V}_k| - \frac{1}{2} \text{tr}(\Psi_0 \mathbf{V}_k^{-1}) + \sum_{j=1}^p \gamma_{jk} \left(-\frac{1}{2} \begin{bmatrix} (\beta_m)_j \\ (\alpha_a)_j \end{bmatrix}^T \mathbf{V}_k^{-1} \begin{bmatrix} (\beta_m)_j \\ (\alpha_a)_j \end{bmatrix} \right)$$

$$p(\mathbf{V}_k | \cdot) \sim \text{Inv-Wishart}(\Psi_0 + \sum_{j=1}^p \gamma_{jk} \begin{bmatrix} (\beta_m)_j \\ (\alpha_a)_j \end{bmatrix} \begin{bmatrix} (\beta_m)_j \\ (\alpha_a)_j \end{bmatrix}^T, \sum_{j=1}^p \gamma_{jk} + \nu)$$

Sampling β_a

$$\log p(\beta_a | \cdot) \propto -\frac{\beta_a^2}{2\sigma_a^2} - \sum_{i=1}^n \left\{ \frac{(A_i \beta_a)^2}{2\sigma_e^2} - \sigma_1^{-2} A_i (Y_i - \mathbf{M}_i^T \beta_m - \mathbf{C}_i^T \beta_c) \beta_a \right\}$$

$$p(\beta_a | \cdot) \sim N\left(\frac{\sum_{i=1}^n A_i (Y_i - \mathbf{M}_i^T \beta_m - \mathbf{C}_i^T \beta_c)}{\sigma_e^2 / \sigma_a^2 + \sum_{i=1}^n A_i^2}, \frac{1}{1/\sigma_a^2 + \sum_{i=1}^n A_i^2 / \sigma_e^2} \right)$$

Sampling σ_a^2

$$\log p(\sigma_a^2 | \cdot) \propto -\left(\frac{1}{2} + h_a + 1\right) \log(\sigma_a^2) - \left(\frac{\beta_a^2}{2} + l_a\right) \sigma_a^{-2}$$

$$p(\sigma_a^2 | \cdot) \sim \text{inverse-gamma}\left(\frac{1}{2} + h_a, \frac{\beta_a^2}{2} + l_a\right)$$

Sampling σ_e^2

$$\log p(\sigma_e^2 | \cdot) = -\left(\frac{n}{2} + h_1 + 1\right) \log(\sigma_e^2) - \left(\frac{\sum_{i=1}^n (Y_i - \mathbf{M}_i^T \beta_m - A_i \beta_a - \mathbf{C}_i^T \beta_c)^2}{2} + l_1\right) \sigma_e^{-2}$$

$$p(\sigma_e^2 | \cdot) \sim \text{inverse-gamma}\left(\frac{n}{2} + h_1, \frac{\sum_{i=1}^n (Y_i - \mathbf{M}_i^T \beta_m - A_i \beta_a - \mathbf{C}_i^T \beta_c)^2}{2} + l_1\right)$$

Sampling σ_g^2

$$\log p(\sigma_g^2 | \cdot) = -\left(\frac{pn}{2} + h_2 + 1\right) \log(\sigma_g^2) - \left(\frac{\sum_{i=1}^n (\mathbf{M}_i^T - A_i \alpha_a - \mathbf{C}_i^T \alpha_c)(\mathbf{M}_i^T - A_i \alpha_a - \mathbf{C}_i^T \alpha_c)^T}{2} + l_2\right) \sigma_g^{-2}$$

$$p(\sigma_g^2|\cdot) \sim \text{inverse-gamma}\left(\frac{pn}{2} + h_2, \frac{\sum_{i=1}^n (\mathbf{M}_i^T - A_i \boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)(\mathbf{M}_i^T - A_i \boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)^T}{2} + l_2\right)$$

Sampling β_{cw}

$$\begin{aligned} \log p(\beta_{cw}|\cdot) &= - \sum_{i=1}^n \left\{ \frac{(C_{iw} \beta_{cw})^2}{2\sigma_e^2} + \sigma_e^{-2} C_{iw} (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i \beta_a - \sum_{s \neq w} C_{is} \beta_{cs}) \beta_{cw} \right\} \\ p(\beta_{cw}|\cdot) &= N\left(\frac{\sum_{i=1}^n C_{iw} (Y_i - A_i \beta_a - \mathbf{M}_i^T \boldsymbol{\beta}_m - \sum_{s \neq w} C_{is} \beta_{cs})}{\sum_{i=1}^n C_{iw}^2}, \frac{\sigma_e^2}{\sum_{i=1}^n C_{iw}^2}\right) \end{aligned}$$

Sampling $(\boldsymbol{\alpha}_{cw})_j$

$$\begin{aligned} \log p((\boldsymbol{\alpha}_{cw})_j|\cdot) &= - \sum_{i=1}^n \left\{ \frac{(C_{iw} (\boldsymbol{\alpha}_{cw})_j)^2}{2\sigma_g^2} + \sigma_g^{-2} C_{iw} (M_i^{(j)} - A_i \alpha_{aj} - \sum_{s \neq w} C_{is} (\boldsymbol{\alpha}_{cs})_j) (\boldsymbol{\alpha}_{cw})_j \right\} \\ p((\boldsymbol{\alpha}_{cw})_j|\cdot) &= N\left(\frac{\sum_{i=1}^n C_{iw} (M_i^{(j)} - A_i \alpha_{aj} - \sum_{s \neq w} C_{is} (\boldsymbol{\alpha}_{cs})_j)}{\sum_{i=1}^n C_{iw}^2}, \frac{\sigma_g^2}{\sum_{i=1}^n C_{iw}^2}\right) \end{aligned}$$

3 Posterior Sampling Algorithm Details for Product Threshold Gaussian (PTG)

Prior

Let $\boldsymbol{\Theta}_{PTG} = (\boldsymbol{\beta}_m, \boldsymbol{\alpha}_a, \tilde{\boldsymbol{\beta}}_m, \tilde{\boldsymbol{\alpha}}_a, \tau_\beta^2, \tau_\alpha^2, \beta_a, \boldsymbol{\beta}_c, \boldsymbol{\alpha}_c, \sigma_e^2, \boldsymbol{\Sigma})$ denote all the unknown parameters in the model. Under the PTG prior, the joint log posterior distribution is,

$$\begin{aligned} \log P(\boldsymbol{\Theta}_{PTG} | \{Y_i, \mathbf{M}_i, A_i, \mathbf{C}_i\}_{i=1}^n) &\propto \sum_{i=1}^n \log P(Y_i | \mathbf{M}_i, \boldsymbol{\beta}_m, \sigma_e^2, \beta_a, \boldsymbol{\beta}_c, A_i, \mathbf{C}_i) + \log P(\mathbf{M}_i | \boldsymbol{\alpha}_a, \boldsymbol{\alpha}_c, \boldsymbol{\Sigma}, A_i, \mathbf{C}_i) \\ &\quad + \log P(\boldsymbol{\Theta}_{PTG}) \\ &= \sum_{i=1}^n -\frac{1}{2} \log \sigma_e^2 - \frac{1}{2\sigma_e^2} (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i \beta_a - \mathbf{C}_i^T \boldsymbol{\beta}_c)^2 \\ &\quad - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i)^T \boldsymbol{\Sigma}^{-1} (\mathbf{M}_i - A_i \boldsymbol{\alpha}_a - \boldsymbol{\alpha}_c \mathbf{C}_i) \\ &\quad + \sum_{i=1}^p -\frac{1}{2} \log \tau_\beta^2 - \frac{(\tilde{\boldsymbol{\beta}}_m)_j^2}{2\tau_\beta^2} + \sum_{i=1}^p -\frac{1}{2} \log \tau_\alpha^2 - \frac{(\tilde{\boldsymbol{\alpha}}_a)_j^2}{2\tau_\alpha^2} \\ &\quad - \frac{q}{2} \log 2\pi \sigma_c^2 - \frac{\boldsymbol{\beta}_c^T \boldsymbol{\beta}_c}{2\sigma_c^2} - \frac{pq}{2} \log 2\pi \sigma_c^2 - \sum_{j=1}^p \frac{\boldsymbol{\alpha}_{cj}^T \boldsymbol{\alpha}_{cj}}{2\sigma_c^2} \end{aligned}$$

Sampling $(\boldsymbol{\beta}_m)_j$

For $(\tilde{\beta}_m)_j$, we denote its threshold conditional on the other parameters as

$$u_{(\tilde{\beta}_m)_j} = \begin{cases} \min(\lambda_1, \lambda_0/|(\tilde{\alpha}_a)_j|), & \text{for } (\tilde{\alpha}_a)_j \neq 0 \\ \lambda_1, & \text{for } (\tilde{\alpha}_a)_j = 0 \end{cases}$$

$$\log p((\tilde{\beta}_m)_j | |(\tilde{\beta}_m)_j| < u_{(\tilde{\beta}_m)_j}) \propto -(\tilde{\beta}_m)_j^2 / (2\tau_\beta^2)$$

$$(\tilde{\beta}_m)_j | |(\tilde{\beta}_m)_j| < u_{(\tilde{\beta}_m)_j} \sim TN(0, \tau_\beta^2, -u_{(\tilde{\beta}_m)_j}, u_{(\tilde{\beta}_m)_j})$$

where $TN(\mu, \sigma^2, a, b)$ denotes a truncated normal distribution with mean μ , variance σ^2 truncated between $[a, b]$.

$$\begin{aligned} & \log p((\tilde{\beta}_m)_j | |(\tilde{\beta}_m)_j| \geq u_{(\tilde{\beta}_m)_j}) \\ & \propto -\frac{(\tilde{\beta}_m)_j^2}{2\tau_\beta^2} - \sum_{i=1}^n \left\{ \frac{(M_i^{(j)}(\tilde{\beta}_m)_j)^2}{2\sigma_e^2} + \sigma_e^{-2} M_i^{(j)} (Y_i - A_i\beta_a - \sum_{s \neq j} M_i^{(s)}(\tilde{\beta}_m)_s - \mathbf{C}_i^T \beta_c) (\tilde{\beta}_m)_j \right\} \end{aligned}$$

$$(\tilde{\beta}_m)_j | (\tilde{\beta}_m)_j \geq u_{(\tilde{\beta}_m)_j} \sim TN(\mu_{mj}, s_{mj}^2, u_{(\tilde{\beta}_m)_j}, \infty)$$

$$(\tilde{\beta}_m)_j | (\tilde{\beta}_m)_j \leq -u_{(\tilde{\beta}_m)_j} \sim TN(\mu_{mj}, s_{mj}^2, -\infty, -u_{(\tilde{\beta}_m)_j})$$

$$\mu_{mj} = \frac{\sum_{i=1}^n M_i^{(j)} (Y_i - A_i\beta_a - \sum_{s \neq j} M_i^{(s)}(\tilde{\beta}_m)_s - \mathbf{C}_i^T \beta_c)}{\sigma_e^2 / \tau_\beta^2 + \sum_{i=1}^n (M_i^{(j)})^2}, \quad s_{mj}^2 = \frac{1}{1/\tau_\beta^2 + \sum_{i=1}^n (M_i^{(j)})^2 / \sigma_e^2}$$

And,

$$\begin{aligned} p(|(\tilde{\beta}_m)_j| < u_{(\tilde{\beta}_m)_j}) &= \frac{B_1}{B_1 + B_2 + B_3} \\ p((\tilde{\beta}_m)_j \geq u_{(\tilde{\beta}_m)_j}) &= \frac{B_2}{B_1 + B_2 + B_3} \\ p((\tilde{\beta}_m)_j \leq -u_{(\tilde{\beta}_m)_j}) &= \frac{B_3}{B_1 + B_2 + B_3} \end{aligned}$$

where $B_1 = \int_{-u_{(\tilde{\beta}_m)_j}}^{u_{(\tilde{\beta}_m)_j}} \frac{1}{\sqrt{2\pi\tau_\beta^2}} \exp(-\frac{(\tilde{\beta}_m)_j^2}{2\tau_\beta^2}) = 1 - 2\Phi(-\frac{u_{(\tilde{\beta}_m)_j}}{\tau_\beta})$, $\Phi(x)$ is the CDF for standard normal distribution, $B_2 = \exp(\mu_{mj}^2 / 2s_{mj}^2 + \log(s_{mj}) - \log(\tau_\beta))(1 - \Phi(\frac{u_{(\tilde{\beta}_m)_j}}{\tau_\beta}))$, $B_3 =$

$$\exp(\mu_{mj}^2/2s_{mj}^2 + \log(s_{mj}) - \log(\tau_\beta))\Phi(-\frac{u_{(\tilde{\beta}_m)_j}}{\tau_\beta^2}).$$

$$(\beta_m)_j = \begin{cases} (\tilde{\beta}_m)_j, & \text{for } |(\tilde{\beta}_m)_j| \geq u_{(\tilde{\beta}_m)_j} \\ 0, & \text{for } |(\tilde{\beta}_m)_j| < u_{(\tilde{\beta}_m)_j} \end{cases}$$

Sampling $(\alpha_a)_j$

For $(\tilde{\alpha}_a)_j$, we denote its threshold conditional on the other parameters as

$$u_{(\tilde{\alpha}_a)_j} = \begin{cases} \min(\lambda_2, \lambda_0/|(\tilde{\beta}_m)_j|), & \text{for } (\tilde{\beta}_m)_j \neq 0 \\ \lambda_2, & \text{for } (\tilde{\beta}_m)_j = 0 \end{cases}$$

$$\log p((\tilde{\alpha}_a)_j | |(\tilde{\alpha}_a)_j| < u_{(\tilde{\alpha}_a)_j}) \propto -(\tilde{\alpha}_a)_j^2/(2\tau_\alpha^2)$$

$$(\tilde{\alpha}_a)_j | |(\tilde{\alpha}_a)_j| < u_{(\tilde{\alpha}_a)_j} \sim TN(0, \tau_\alpha^2, -u_{(\tilde{\alpha}_a)_j}, u_{(\tilde{\alpha}_a)_j})$$

$$\log p((\tilde{\alpha}_a)_j | |(\tilde{\alpha}_a)_j| \geq u_{(\tilde{\alpha}_a)_j}) \propto -\frac{(\tilde{\alpha}_a)_j^2}{2\tau_\alpha^2} - \sum_{i=1}^n \left\{ \frac{(A_i(\tilde{\alpha}_a)_j)^2}{2\sigma_g^2} + \sigma_g^{-2} A_i (M_i^{(j)} - (\alpha_c \mathbf{C}_i)_j) (\tilde{\alpha}_a)_j \right\}$$

$$(\tilde{\alpha}_a)_j | (\tilde{\alpha}_a)_j \geq u_{(\tilde{\alpha}_a)_j} \sim TN(\mu_{aj}, s_{aj}^2, u_{(\tilde{\alpha}_a)_j}, \infty)$$

$$(\tilde{\alpha}_a)_j | (\tilde{\alpha}_a)_j \leq -u_{(\tilde{\alpha}_a)_j} \sim TN(\mu_{aj}, s_{aj}^2, -\infty, -u_{(\tilde{\alpha}_a)_j})$$

$$\mu_{aj} = \frac{\sum_{i=1}^n A_i (M_i^{(j)} - (\alpha_c \mathbf{C}_i)_j)}{\sigma_g^2/\tau_\alpha^2 + \sum_{i=1}^n A_i^2}, s_{aj}^2 = \frac{1}{1/\tau_\alpha^2 + \sum_{i=1}^n A_i^2/\sigma_g^2}$$

And,

$$\begin{aligned} p(|(\tilde{\alpha}_a)_j| < u_{(\tilde{\alpha}_a)_j}) &= \frac{A_1}{A_1 + A_2 + A_3} \\ p((\tilde{\alpha}_a)_j \geq u_{(\tilde{\alpha}_a)_j}) &= \frac{A_2}{A_1 + A_2 + A_3} \\ p((\tilde{\alpha}_a)_j \leq -u_{(\tilde{\alpha}_a)_j}) &= \frac{A_3}{A_1 + A_2 + A_3} \end{aligned}$$

where $A_1 = \int_{-u_{(\tilde{\alpha}_a)_j}}^{u_{(\tilde{\alpha}_a)_j}} \frac{1}{\sqrt{2\pi\tau_\alpha^2}} \exp(-\frac{(\tilde{\alpha}_a)_j^2}{2\tau_\alpha^2}) = 1 - 2\Phi(-\frac{u_{(\tilde{\alpha}_a)_j}}{\tau_\alpha^2})$, $\Phi(x)$ is the CDF for standard normal distribution, $A_2 = \exp(\mu_{aj}^2/2s_{aj}^2 + \log(s_{aj}) - \log(\tau_\alpha))(1 - \Phi(\frac{u_{(\tilde{\alpha}_a)_j}}{\tau_\alpha^2}))$, $A_3 =$

$$\exp(\mu_{\alpha_j}^2/2s_{\alpha_j}^2 + \log(s_{\alpha_j}) - \log(\tau_\alpha))\Phi(-\frac{u(\tilde{\alpha}_a)_j}{\tau_\alpha^2}).$$

$$(\alpha_a)_j = \begin{cases} (\tilde{\alpha}_a)_j, & \text{for } |(\tilde{\alpha}_a)_j| \geq u(\tilde{\alpha}_a)_j \\ 0, & \text{for } |(\tilde{\alpha}_a)_j| < u(\tilde{\alpha}_a)_j \end{cases}$$

Sampling β_a

$$\log p(\beta_a|\cdot) \propto -\frac{\beta_a^2}{2\sigma_a^2} - \sum_{i=1}^n \left\{ \frac{(A_i\beta_a)^2}{2\sigma_1^2} - \sigma_1^{-2} A_i (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - \mathbf{C}_i^T \boldsymbol{\beta}_c) \beta_a \right\}$$

$$p(\beta_a|\cdot) \sim N\left(\frac{\sum_{i=1}^n A_i (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - \mathbf{C}_i^T \boldsymbol{\beta}_c)}{\sigma_1^2/\sigma_a^2 + \sum_{i=1}^n A_i^2}, \frac{1}{1/\sigma_a^2 + \sum_{i=1}^n A_i^2/\sigma_1^2}\right)$$

Sampling σ_a^2

$$\log p(\sigma_a^2|\cdot) \propto -\left(\frac{1}{2} + h_a + 1\right)\log(\sigma_a^2) - \left(\frac{\beta_a^2}{2} + l_a\right)\sigma_a^{-2}$$

$$p(\sigma_a^2|\cdot) \sim \text{inverse-gamma}\left(\frac{1}{2} + h_a, \frac{\beta_a^2}{2} + l_a\right)$$

Sampling σ_e^2

$$\log p(\sigma_e^2|\cdot) = -\left(\frac{n}{2} + h_1 + 1\right)\log(\sigma_e^2) - \left(\frac{\sum_{i=1}^n (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i\beta_a - \mathbf{C}_i^T \boldsymbol{\beta}_c)^2}{2} + l_1\right)\sigma_e^{-2}$$

$$p(\sigma_e^2|\cdot) \sim \text{inverse-gamma}\left(\frac{n}{2} + h_1, \frac{\sum_{i=1}^n (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i\beta_a - \mathbf{C}_i^T \boldsymbol{\beta}_c)^2}{2} + l_1\right)$$

Sampling σ_g^2

$$\log p(\sigma_g^2|\cdot) = -\left(\frac{pn}{2} + h_2 + 1\right)\log(\sigma_g^2) - \left(\frac{\sum_{i=1}^n (\mathbf{M}_i^T - A_i\boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)(\mathbf{M}_i^T - A_i\boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)^T}{2} + l_2\right)\sigma_g^{-2}$$

$$p(\sigma_g^2|\cdot) \sim \text{inverse-gamma}\left(\frac{pn}{2} + h_2, \frac{\sum_{i=1}^n (\mathbf{M}_i^T - A_i\boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)(\mathbf{M}_i^T - A_i\boldsymbol{\alpha}_a - \mathbf{C}_i^T \boldsymbol{\alpha}_c)^T}{2} + l_2\right)$$

Sampling τ_β^2

$$\log p(\tau_\beta^2|\cdot) = -\left(\frac{q}{2} + k_m + 1\right)\log(\tau_\beta^2) - \left(\frac{\sum_{j=1}^q (\tilde{\boldsymbol{\beta}}_m)_j^2}{2} + l_m\right)\tau_\beta^{-2}$$

$$p(\tau_\beta^2|\cdot) \sim \text{inverse-gamma}\left(\frac{q}{2} + k_m, \frac{\sum_{j=1}^q (\tilde{\beta}_m)_j^2}{2} + l_m\right)$$

Sampling τ_α^2

$$\begin{aligned} \log p(\tau_\alpha^2|\cdot) &= -\left(\frac{q}{2} + k_{ma} + 1\right)\log(\tau_\alpha^2) - \left(\frac{\sum_{j=1}^q (\tilde{\alpha}_a)_j^2}{2} + l_{ma}\right)\tau_\alpha^{-2} \\ p(\tau_\alpha^2|\cdot) &\sim \text{inverse-gamma}\left(\frac{q}{2} + k_{ma}, \frac{\sum_{j=1}^q (\tilde{\alpha}_a)_j^2}{2} + l_{ma}\right) \end{aligned}$$

Sampling β_{cw}

$$\begin{aligned} \log p(\beta_{cw}|\cdot) &= -\sum_{i=1}^n \left\{ \frac{(C_{iw}\beta_{cw})^2}{2\sigma_e^2} + \sigma_e^{-2} C_{iw} (Y_i - \mathbf{M}_i^T \boldsymbol{\beta}_m - A_i \beta_a - \sum_{s \neq w} C_{is} \beta_{cs}) \beta_{cw} \right\} \\ p(\beta_{cw}|\cdot) &= N\left(\frac{\sum_{i=1}^n C_{iw} (Y_i - A_i \beta_a - \mathbf{M}_i^T \boldsymbol{\beta}_m - \sum_{s \neq w} C_{is} \beta_{cs})}{\sum_{i=1}^n C_{iw}^2}, \frac{\sigma_e^2}{\sum_{i=1}^n C_{iw}^2}\right) \end{aligned}$$

Sampling $(\boldsymbol{\alpha}_{cw})_j$

$$\begin{aligned} \log p((\boldsymbol{\alpha}_{cw})_j|\cdot) &= -\sum_{i=1}^n \left\{ \frac{(C_{iw}(\boldsymbol{\alpha}_{cw})_j)^2}{2\sigma_g^2} + \sigma_g^{-2} C_{iw} (M_i^{(j)} - A_i \alpha_{aj} - \sum_{s \neq w} C_{is} (\boldsymbol{\alpha}_{cs})_j) (\boldsymbol{\alpha}_{cw})_j \right\} \\ p((\boldsymbol{\alpha}_{cw})_j|\cdot) &= N\left(\frac{\sum_{i=1}^n C_{iw} (M_i^{(j)} - A_i \alpha_{aj} - \sum_{s \neq w} C_{is} (\boldsymbol{\alpha}_{cs})_j)}{\sum_{i=1}^n C_{iw}^2}, \frac{\sigma_g^2}{\sum_{i=1}^n C_{iw}^2}\right) \end{aligned}$$

4 Effects Distribution and Additional Results in Simulations

Effects Distribution

To better understand the generated effects under the three different data generating mechanism in the simulation Setting (A)-(C), we examine the corresponding distributions of the simulated non-zero marginal effects, $(\boldsymbol{\beta}_m)_j$ (or $(\boldsymbol{\alpha}_a)_j$) and indirect effects, $(\boldsymbol{\beta}_m)_j(\boldsymbol{\alpha}_a)_j$ in Figure S1.

The PTG prior model essentially produces effects truncated away from zero (Setting (A)), where the thresholding parameter $\lambda = (\lambda_0, \lambda_1, \lambda_2)$ is determined by the proportion of non-zero effects. For example, choosing $\lambda_0 = |\tilde{\alpha}_a \tilde{\beta}_m|^{(95)}$, $\lambda_1 = |\tilde{\beta}_m|^{(85)}$, $\lambda_2 = |\tilde{\alpha}_a|^{(93)}$ approximately makes $\pi_1 = 0.05$, $\pi_2 = 0.10$, $\pi_3 = 0.05$, $\pi_4 = 0.80$. The relatively small non-zero

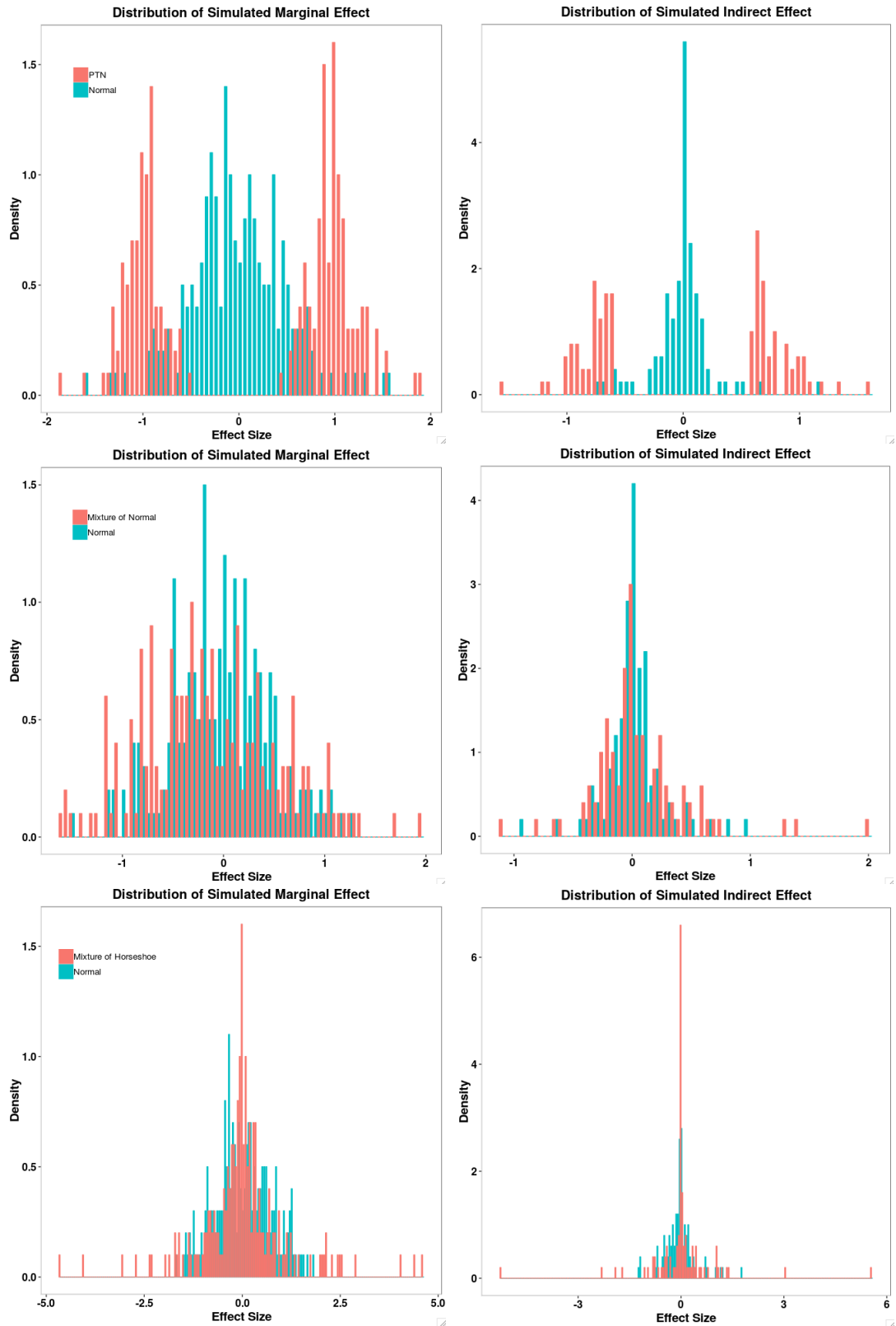


Figure S1: The distributions of the simulated non-zero marginal effects, $(\beta_m)_j$ (or $(\alpha_a)_j$) and indirect effects, $(\beta_m)_j(\alpha_a)_j$ under the three simulation settings when $n = 100, p = 200$. Each row represents one scenario, i.e. effects under prior model PTG, GMM and Mixture of Horseshoe. We include marginal effects from normals with the same variances as the simulation distributions and the corresponding indirect effects as a comparison.

marginal effects are picked up by its indirect effects exceeding the product threshold. The Setting (B) with four components of bivariate Gaussian mixture is straightforward, and the resulting indirect effects distribute as a product of two normal distributions. Under the Setting (C), we can see that the horseshoe distribution has a tall spike near zero and heavy tails on large effects, and this generates uneven effects different from either PTG or GMM prior model. The distribution of the corresponding indirect effects show a stronger contrast between small and large effects.

Empirical FDR Results

As a practical procedure, we suggest a cutoff on the posterior inclusion probabilities (PIP) to identify a significance threshold for declaring active mediators. To evaluate the performance of this significance rule, we report the empirical FDR and TPR in Table S7 and S2 under all the simulation scenarios. We find that at $PIP = 0.5$ cutoff, the two proposed methods, PTG and GMM, exhibit good selection performance while maintaining a reasonable FDR in most scenarios. At $PIP = 0.9$ cutoff, the two methods provide over conservative estimates of FDR, leading to reduced power in mediator selection. Therefore, we will use the 0.5 cutoff on the PIPs as a selection criterion in the following applications.

Computing Time

We performed simulations on a single core of Intel(R) Xeon(R) Platinum 8176 CPU @ 2.10GHz, and the runtime comparison of the proposed methods is shown in Table S3. For both the small sample scenario with $n = 100$, $p = 200$, and the large sample scenario with $n = 1000$, $p = 2000$, the proposed algorithms can be finished in a reasonable amount of time. We still acknowledge that future development of new algorithms and/or new methods will likely be required to scale our method to handle thousands of individuals and millions of mediators.

Estimation Bias

Besides MSE, bias is another important metric for effect estimation. We included here separate result tables for the bias metric on joint NIE, NDE and TE for all the simulation

Method	TPR(FDR=0.1)	TPR(PIP>0.9)	FDR(PIP>0.9)	TPR(PIP>0.5)	FDR(PIP>0.5)
<i>n = 100, p = 200, p₁₁ = 10, fixed effects (I)</i>					
PTG	0.54(0.025)	0.27(0.017)	0.03(0.014)	0.55(0.017)	0.13(0.016)
GMM	0.42(0.023)	0.17(0.022)	0.03(0.021)	0.44(0.023)	0.16(0.017)
<i>n = 100, p = 200, p₁₁ = 10, fixed effects (II)</i>					
PTG	0.34(0.017)	0.27(0.008)	0.04(0.019)	0.37(0.013)	0.14(0.019)
GMM	0.39(0.020)	0.21(0.010)	0.03(0.016)	0.39(0.013)	0.11(0.017)
<i>n = 100, p = 200, p₁₁ = 10, PTG</i>					
PTG	0.45(0.020)	0.19(0.014)	0.01(0.007)	0.49(0.018)	0.18(0.015)
GMM	0.43(0.015)	0.26(0.011)	0.03(0.012)	0.45(0.014)	0.11(0.012)
<i>n = 100, p = 200, p₁₁ = 10, Gaussian</i>					
PTG	0.38(0.008)	0.26(0.008)	0.01(0.006)	0.56(0.010)	0.39(0.011)
GMM	0.41(0.006)	0.27(0.005)	0.01(0.002)	0.35(0.006)	0.06(0.008)
<i>n = 100, p = 200, p₁₁ = 10, Horseshoe</i>					
PTG	0.30(0.015)	0.24(0.014)	0.08(0.016)	0.37(0.016)	0.38(0.019)
GMM	0.33(0.011)	0.26(0.011)	0.03(0.008)	0.35(0.012)	0.16(0.014)

Table S1: Empirical estimates of TPR and FDR in simulations under $n = 100, p = 200, p_{11}$ is the number of true active mediators. The results are based on 200 replicates for each setting, and the standard errors are shown within parentheses. TPR(FDR=0.1) is the true positive rate controlled at a fixed FDR of 10%; TPR(PIP>0.9) and FDR(PIP>0.9) are the empirical estimates when the PIP threshold for declaring active mediators is 0.9; TPR(PIP>0.5) and FDR(PIP>0.5) are the empirical estimates when the PIP threshold for declaring active mediators is 0.5.

Method	TPR(FDR=0.1)	TPR(PIP>0.9)	FDR(PIP>0.9)	TPR(PIP>0.5)	FDR(PIP>0.5)
<i>n = 1000, p = 2000, p₁₁ = 100, fixed effects (I)</i>					
PTG	0.64(0.008)	0.49(0.017)	0.01(0.002)	0.55(0.017)	0.06(0.016)
GMM	0.61(0.009)	0.40(0.004)	0.01(0.003)	0.55(0.005)	0.07(0.010)
<i>n = 1000, p = 2000, p₁₁ = 100, fixed effects (II)</i>					
PTG	0.40(0.008)	0.20(0.004)	0.01(0.003)	0.37(0.012)	0.07(0.010)
GMM	0.48(0.006)	0.29(0.003)	0.01(0.002)	0.43(0.004)	0.06(0.007)
<i>n = 1000, p = 2000, p₁₁ = 100, PTG</i>					
PTG	0.40(0.008)	0.19(0.004)	0.01(0.011)	0.44(0.007)	0.13(0.006)
GMM	0.37(0.010)	0.10(0.004)	0.05(0.008)	0.47(0.006)	0.17(0.007)
<i>n = 1000, p = 2000, p₁₁ = 100, Gaussian</i>					
PTG	0.42(0.006)	0.20(0.005)	0.03(0.002)	0.51(0.005)	0.17(0.004)
GMM	0.51(0.007)	0.36(0.005)	0.01(0.002)	0.49(0.006)	0.10(0.004)
<i>n = 1000, p = 2000, p₁₁ = 100, Horseshoe</i>					
PTG	0.29(0.008)	0.30(0.004)	0.05(0.006)	0.39(0.008)	0.24(0.004)
GMM	0.38(0.007)	0.35(0.004)	0.03(0.003)	0.45(0.004)	0.18(0.015)

Table S2: Empirical estimates of TPR and FDR in simulations under $n = 1000, p = 2000, p_{11}$ is the number of true active mediators. The results are based on 200 replicates for each setting, and the standard errors are shown within parentheses. TPR(FDR=0.1) is the true positive rate controlled at a fixed FDR of 10%; TPR(PIP>0.9) and FDR(PIP>0.9) are the empirical estimates when the PIP threshold for declaring active mediators is 0.9; TPR(PIP>0.5) and FDR(PIP>0.5) are the empirical estimates when the PIP threshold for declaring active mediators is 0.5.

settings discussed in the main paper (see Table S4, S5), and summarized the results in the simulation results section of the main paper.

Method	$n = 100, p = 200$	$n = 1000, p = 2000$
PTG	30.5sec	23.0min
GMM	88.8sec	29.8min

Table S3: The average runtime of the proposed methods for $n = 100, p = 200$ and $n = 1000, p = 2000$ in the simulations. Comparison was carried out on a single core of Intel(R) Xeon(R) Platinum 8176 CPU @ 2.10GHz. For the proposed methods, we in total ran 150,000 iterations.

Table S4: Estimation bias for fixed effect simulations under $n = 100, p = 200$ and $n = 1000, p = 2000, p_{11}$ is the number of truly active mediators. The results are based on 200 replicates for each setting.

$n = 100, p = 200, p_{11} = 10, fixed\ effects\ (I)$				$fixed\ effects\ (II)$		
<i>Method</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>
PTG (0.15,0.4,0.4)	-0.33	0.21	-0.12	1.34	-1.15	0.19
GMM	0.15	-0.19	-0.04	0.92	-0.84	0.08
BAMA	-0.17	0.10	-0.07	1.12	-1.01	0.11
Bi-BLasso	-0.41	0.39	-0.02	0.94	-0.83	0.11
PathLasso	-0.34	0.23	-0.11	1.61	-1.34	0.27
Bi-Lasso	-0.35	0.25	-0.10	1.41	-1.17	0.24
HIMA	-0.20	-0.01	-0.21	1.24	-0.92	0.32
Univariate	-12.01	11.52	-0.49	-10.23	9.82	-0.41
$n = 1000, p = 2000, p_{11} = 100, fixed\ effects\ (I)$				$fixed\ effects\ (II)$		
<i>Method</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>
PTG (0.15,0.4,0.4)	-0.19	0.15	-0.04	0.10	-0.05	0.05
GMM	0.06	-0.08	-0.02	0.01	0.03	0.04
BAMA	-1.15	1.14	-0.01	0.09	-0.03	0.06
Bi-BLasso	0.77	-0.74	0.03	0.31	-0.34	-0.03
PathLasso	-2.40	2.47	0.07	0.42	-0.46	-0.04
Bi-Lasso	-2.23	2.27	0.04	0.23	-0.27	-0.04
HIMA	-3.27	2.52	-0.75	0.54	0.58	1.12
Univariate	-16.21	18.73	2.52	-7.21	6.93	-0.28

Table S5: Estimation bias for random effect simulations under $n = 100, p = 200$ and $n = 1000, p = 2000$, p_{11} is the number of truly active mediators. The results are based on 200 replicates for each setting.

$n = 100, p = 200, p_{11} = 10, PTG, \sigma_u^2 = 0.3$				$Gaussian, \sigma^2 = 0.3$			$Horseshoe, \sigma^2 = 0.5, b = 3$		
<i>Method</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>
PTG	-0.086	0.079	-0.007	-0.403	0.355	-0.048	-0.190	0.103	-0.087
GMM	0.287	-0.292	-0.005	0.319	-0.322	-0.003	0.188	-0.203	-0.015
BAMA	0.187	-0.194	-0.007	-0.194	0.120	-0.074	0.114	-0.147	-0.033
Bi-BLasso	-0.233	0.261	0.028	-0.255	0.215	-0.040	-0.540	0.512	-0.028
PathLasso	-0.284	0.252	-0.032	-0.230	0.154	-0.076	-0.692	0.619	-0.073
Bi-Lasso	-0.298	0.307	0.009	-0.076	0.042	-0.034	-0.999	0.953	-0.046
HIMA	-0.129	0.282	0.153	-0.168	-0.169	-0.337	-0.741	-0.902	-1.643
Univariate	-12.029	11.824	-0.205	-7.940	7.794	-0.146	-14.108	13.872	-0.236
$n = 1000, p = 2000, p_{11} = 100, PTG, \sigma_u^2 = 0.1$				$Gaussian, \sigma^2 = 0.1$			$Horseshoe, \sigma^2 = 0.3, b = 3$		
<i>Method</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>	<i>NIE</i>	<i>NDE</i>	<i>TE</i>
PTG	0.026	-0.009	0.017	-0.780	0.708	-0.072	-0.163	0.195	0.032
GMM	0.020	-0.013	0.007	0.049	-0.047	0.002	0.088	-0.117	-0.029
BAMA	0.046	-0.057	-0.011	-0.287	0.263	-0.024	-0.204	0.139	-0.065
Bi-BLasso	-0.067	-0.027	-0.094	-0.656	0.614	-0.042	-0.277	0.098	-0.179
PathLasso	0.390	-0.428	-0.038	0.410	-0.383	0.027	-7.091	7.049	-0.042
Bi-Lasso	0.004	0.009	0.013	-0.763	0.795	0.032	4.650	-4.693	-0.043
HIMA	0.033	0.281	0.314	-1.438	-0.300	-1.738	-46.334	76.720	30.386
Univariate	-6.240	6.179	-0.061	-11.124	11.099	-0.025	-17.095	17.046	-0.049

Data-Adaptive Uniform Priors on λ 's

Alternative to fixing the threshold parameter λ , we also consider a data-adaptive uniform prior that favors large positive value on λ 's. From our experience, the lower bound of such informative prior should be far away from zero, otherwise the method tends to include many false positives. Therefore, we first fit the Lasso method and then use the posterior quantiles (e.g. 95% to 99%) of the estimated $|\beta_m|$, $|\alpha_a|$ to determine the range of corresponding λ 's. To be specific, *a priori*, $\lambda_1 \sim U[|\hat{\beta}_m|^{(95\%)}, |\hat{\beta}_m|^{(99\%)}]$, $\lambda_2 \sim U[|\hat{\alpha}_a|^{(95\%)}, |\hat{\alpha}_a|^{(99\%)}]$, and we always set the value of λ_0 as $\lambda_1 \lambda_2$. In Figure S2, we visualize the joint distribution of β_{mj} and α_{aj} under three different prior choices for λ_1 and λ_2 : (1) λ_1, λ_2 are fixed to be *a priori* determined values; (2) λ_1, λ_2 are specified to follow uniform distributions, with the lower thresholds set to be zero and the upper thresholds set to values corresponding to 0.1 prior inclusion probability; or (3) λ_1, λ_2 are specified to follow uniform distributions, with both the lower and upper thresholds determined in a data-driven fashion as explained in the paragraph. We refer to the first option as the fixed value option, the second option as uniform prior option, and the third option as uniform prior with data-adaptive thresholds

option. There is heavy density on the (0,0) in all the plots.

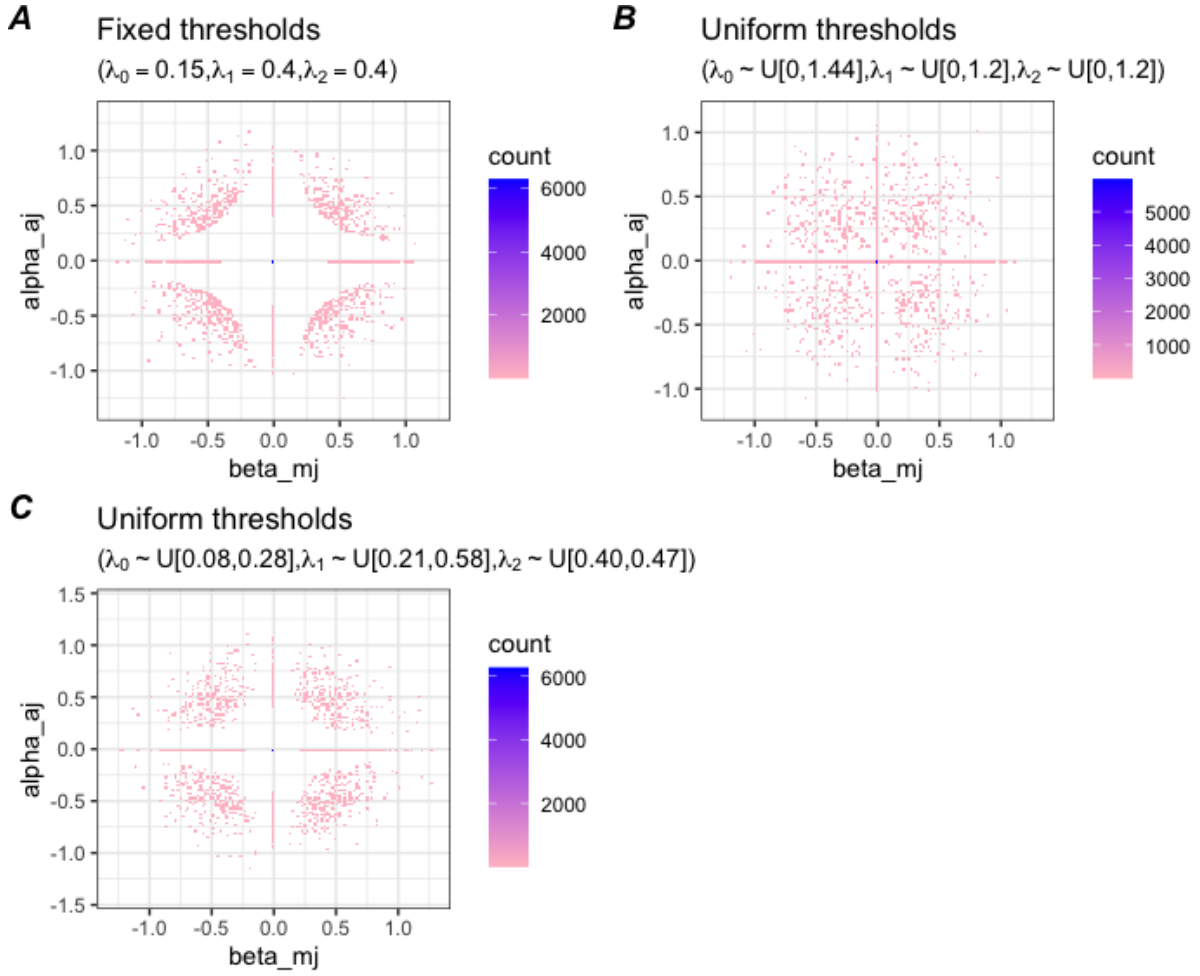


Figure S2: Visualizations of the joint distribution $(\beta_{mj}, \alpha_{aj})$ under fixed value option, uniform prior options after integrating out the λ 's. **A:** under fixed thresholds with $\lambda_0 = 0.15, \lambda_1 = \lambda_2 = 0.4$; **B:** the λ 's are assigned uniform priors with lower bounds being zero, and upper bounds tuned for 0.1 prior inclusion probability; **C:** the λ 's are assigned uniform priors with the range determined by the posterior quantiles (95% to 99%) of the estimated $|\beta_m|, |\alpha_a|$ from Lasso.

We examine the performance of using these three prior choices for λ_1 and λ_2 in the fixed effect simulations. The comparison among the three prior choices is shown in Table S6. The results indicate that uniform priors adequately large lower bounds (e.g. 95% quantiles) on λ_1 (and λ_2) can boost the selection power and estimation accuracy. The two thresholds specified this way also cover a reasonably wide range, e.g. $0.2 \sim 0.6$ when the non-zero true effect is 0.5 in our simulation. With the empirical FDR results shown in Table S7, we find that uniform priors with range of 95% to 99% quantiles tend to produce more conservative results than fixed thresholds.

Table S6: Simulation results for fixed effects under $n = 100, p = 200$ and $n = 1000, p = 2000$, p_{11} is the number of truly active mediators. TPR: true positive rate at false discovery rate (FDR) = 0.10. $MSE_{\text{non-null}}$: mean squared error for the indirect effects of truly active mediators. MSE_{null} : mean squared error for the indirect effects of truly inactive mediators. The results are based on 200 replicates for each setting, and the standard errors are shown within parentheses.

$n = 100, p = 200, p_{11} = 10, \text{fixed effects (I)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-4}$
PTG (fixed values, 0.15,0.4,0.4)	0.99(0.001)	0.52 (0.026)	0.043	0.395
PTG (uniform, 90% to 99%)	0.98(0.001)	0.45 (0.025)	0.046	0.390
PTG (uniform, 95% to 99%)	0.99(0.001)	0.54 (0.027)	0.048	0.246
GMM	0.98(0.001)	0.44 (0.022)	0.047	1.409
BAMA	0.97(0.002)	0.38(0.021)	0.063	2.471
$n = 100, p = 200, p_{11} = 10, \text{fixed effects (II)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-4}$
PTG (fixed values, 0.15,0.4,0.4)	0.96(0.003)	0.35 (0.016)	0.073	0.309
PTG (uniform, 90% to 99%)	0.96(0.003)	0.35 (0.018)	0.077	0.342
PTG (uniform, 95% to 99%)	0.94(0.005)	0.35 (0.016)	0.079	0.153
GMM	0.96(0.003)	0.37 (0.017)	0.062	0.940
BAMA	0.95(0.003)	0.31(0.015)	0.075	2.389
$n = 1000, p = 2000, p_{11} = 100, \text{fixed effects (I)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-4}$
PTG (fixed values, 0.15,0.4,0.4)	0.98(0.001)	0.64 (0.008)	0.028	0.070
PTG (uniform, 90% to 99%)	0.96(0.001)	0.45 (0.010)	0.051	1.206
PTG (uniform, 95% to 99%)	0.98(0.005)	0.65 (0.011)	0.044	0.050
GMM	0.99(0.001)	0.61 (0.009)	0.023	0.134
BAMA	0.98(0.001)	0.54(0.007)	0.040	0.150
$n = 1000, p = 2000, p_{11} = 100, \text{fixed effects (II)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-6}$
PTG (fixed values, 0.15,0.4,0.4)	0.96(0.002)	0.40 (0.008)	0.008	0.164
PTG (uniform, 90% to 99%)	0.96(0.001)	0.31 (0.006)	0.008	0.249
PTG (uniform, 95% to 99%)	0.96(0.001)	0.37 (0.007)	0.008	0.239
GMM	0.97(0.001)	0.48 (0.006)	0.003	3.437
BAMA	0.95(0.001)	0.35(0.005)	0.005	7.485

Table S7: Empirical estimates of TPR and FDR in simulations under $n = 100, p = 200, p_{11}$ is the number of true active mediators. The results are based on 200 replicates for each setting, and the standard errors are shown within parentheses. TPR(FDR=0.1) is the true positive rate controlled at a fixed FDR of 10%; TPR(PIP>0.9) and FDR(PIP>0.9) are the empirical estimates when the PIP threshold for declaring active mediators is 0.9; TPR(PIP>0.5) and FDR(PIP>0.5) are the empirical estimates when the PIP threshold for declaring active mediators is 0.5.

Method	TPR(FDR=0.1)	TPR(PIP>0.9)	FDR(PIP>0.9)	TPR(PIP>0.5)	FDR(PIP>0.5)
<i>n = 100, p = 200, p₁₁ = 10, fixed effects (I)</i>					
PTG (0.15,0.4,0.4)	0.52(0.026)	0.26(0.017)	0.03(0.015)	0.51(0.017)	0.13(0.014)
PTG (uniform, 0.90)	0.45(0.025)	0.17(0.015)	0.04(0.024)	0.48(0.017)	0.16(0.016)
PTG (uniform, 0.95)	0.54(0.027)	0.19(0.018)	0.02(0.012)	0.44(0.021)	0.10(0.015)
<i>n = 100, p = 200, p₁₁ = 10, fixed effects (II)</i>					
PTG (0.15,0.4,0.4)	0.35(0.016)	0.27(0.009)	0.04(0.017)	0.37(0.014)	0.15(0.016)
PTG (uniform, 0.90)	0.35(0.018)	0.19(0.007)	0.03(0.001)	0.38(0.017)	0.20(0.021)
PTG (uniform, 0.95)	0.35(0.016)	0.16(0.006)	0.01(0.001)	0.30(0.013)	0.11(0.018)
<i>n = 1000, p = 2000, p₁₁ = 100, fixed effects (I)</i>					
PTG (0.15,0.4,0.4)	0.64(0.008)	0.49(0.017)	0.01(0.002)	0.55(0.017)	0.06(0.016)
PTG (uniform, 0.90)	0.45(0.010)	0.38(0.020)	0.08(0.007)	0.32(0.005)	0.05(0.006)
PTG (uniform, 0.95)	0.65(0.011)	0.41(0.004)	0.02(0.002)	0.42(0.004)	0.02(0.002)
<i>n = 1000, p = 2000, p₁₁ = 100, fixed effects (II)</i>					
PTG (0.15,0.4,0.4)	0.40(0.008)	0.20(0.004)	0.01(0.003)	0.37(0.012)	0.07(0.010)
PTG (uniform, 0.90)	0.31(0.006)	0.17(0.005)	0.03(0.003)	0.35(0.009)	0.15(0.005)
PTG (uniform, 0.95)	0.37(0.007)	0.19(0.003)	0.05(0.007)	0.55(0.009)	0.27(0.008)

Sensitivity Analysis

We perform sensitivity analysis to examine how robust the posterior inference is regarding mild changes in terms of the prior choices for λ_1 and λ_2 : (1) λ_1, λ_2 are fixed to be *a priori* determined values; (2) λ_1, λ_2 are specified to follow uniform distributions, with both the lower and upper thresholds determined in a data-driven fashion as explained in the paragraph. We refer to the first option as the fixed value option, and the second option as uniform prior option with data-adaptive thresholds option.

We summarize the results in Table S8. In general, the lambda parameters, especially the lower bounds for β_{mj} and α_{aj} , play an important role in PTG's performance. As the lambda parameters vary, the TPR varies in a reasonable range, and is mostly higher than the other methods, while the MSEs are relatively more robust.

5 Detailed Description of MESA Data

MESA is a population-based longitudinal study designed to identify risk factors for the progression of subclinical cardiovascular disease (CVD) [4]. A total of 6,814 non-Hispanic

Table S8: Sensitivity analysis under $n = 100, p = 200$ with fixed effects, p_{11} is the number of truly active mediators. We include the pre-defined fixed thresholds $(\lambda_0, \lambda_1, \lambda_2)$, or the range of uniform priors under each setting. The results are based on 200 replicates for each setting, and the standard errors are shown within parentheses.

$n = 100, p = 200, p_{11} = 10, \text{fixed effects (I)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-4}$
PTG (fixed values, 0.15,0.4,0.4)	0.99(0.001)	0.52(0.026)	0.043	0.395
PTG (fixed values, 0.2,0.4,0.2)	0.96(0.006)	0.47(0.023)	0.048	0.521
PTG (fixed values, 0.25,0.25,0.25)	0.98(0.001)	0.44(0.024)	0.049	0.243
PTG (fixed values, 0.5,0.1,0.5)	0.97(0.002)	0.35(0.019)	0.042	0.863
PTG (uniform, 90% to 99%)	0.98(0.001)	0.45(0.025)	0.046	0.390
PTG (uniform, 95% to 99%)	0.99(0.001)	0.54(0.027)	0.048	0.246
$n = 100, p = 200, p_{11} = 10, \text{fixed effects (II)}$				
<i>Method</i>	<i>AUC</i>	<i>TPR</i>	$MSE_{\text{non-null}}$	$MSE_{\text{null}} \times 10^{-4}$
PTG (fixed values, 0.15,0.4,0.4)	0.96(0.003)	0.35(0.016)	0.073	0.309
PTG (fixed values, 0.2,0.4,0.2)	0.92(0.005)	0.32(0.016)	0.081	0.223
PTG (fixed values, 0.25,0.25,0.25)	0.96(0.003)	0.31(0.015)	0.085	0.187
PTG (fixed values, 0.5,0.1,0.5)	0.93(0.004)	0.27(0.013)	0.082	0.469
PTG (uniform, 90% to 99%)	0.96(0.003)	0.35(0.018)	0.077	0.342
PTG (uniform, 95% to 99%)	0.94(0.005)	0.35(0.016)	0.079	0.153

white, African-American, Hispanic, and Chinese-American women and men aged 45–84 without clinically apparent CVD were recruited between July 2000 and August 2002 from the following 6 regions in the US: Forsyth County, NC; Northern Manhattan and the Bronx, NY; Baltimore City and Baltimore County, MD; St. Paul, MN; Chicago, IL; and Los Angeles County, CA. Each field center recruited from locally available sources, which included lists of residents, lists of dwellings, and telephone exchanges. Neighborhood socioeconomic disadvantage scores for each neighborhood were created based on a principal components analysis of 16 census-tract level variables from the 2000 US Census. These variables reflect dimensions of education, occupation, income and wealth, poverty, employment, and housing. For the neighborhood measures, we use the cumulative average of the measure across all available MESA examinations. The descriptive statistics for the exposure and outcome can be found in Table S9.

In the MESA data, between April 2010 and February 2012 (corresponding to MESA Exam 5), DNAm were assessed on a random subsample of 1,264 non-Hispanic white, African-American, and Hispanic MESA participants aged 55–94 from the Baltimore, Forsyth County, New York, and St. Paul field centers. After excluding respondents with missing data on one or more variables, we had phenotype and DNAm data from purified monocytes on a

total of 1,225 individuals and we focused on this set of individuals for analysis. The detailed description of DNAm data collection, quantitation and data processing procedures can be found in Liu et al [5]. Briefly, the Illumina HumanMethylation450 BeadChip was used to measure DNAm, and bead-level data were summarized in GenomeStudio. Quantile normalization was performed using the *lumi* package with default settings [6]. Quality control (QC) measures included checks for sex and race/ethnicity mismatches and outlier identification by multidimensional scaling plots. Further probe filtering criteria included: “detected” DNAm levels in <90% of MESA samples (detection p -value cut-off = 0.05), existence of a SNP within 10 base pairs of the target CpG site, overlap with a non-unique region, and suggestions by DMRcate [7] (mostly cross-reactive probes). Those procedures leave us 403,713 autosomal probes for analysis.

For computational reasons, we first selected a subset of CpG sites to be used in the final multivariate mediation analysis model. In particular, for each single CpG site in turn, we fit the following linear mixed model to test the marginal association between the CpG site and the exposure variable:

$$M_i = A_i\psi_a + \mathbf{C}_{1i}^T\psi_c + \mathbf{Z}_i^T\psi_u + \epsilon_i, i = 1, \dots, n \quad (1)$$

where A_i represents neighborhood SES value for the i 'th individual and ψ_a is its coefficient; \mathbf{C}_{1i} is a vector of covariates that include age, gender, race/ethnicity, childhood socioeconomic status, adult socioeconomic status and enrichment scores for each of 4 major blood cell types (neutrophils, B cells, T cells and natural killer cells) to account for potential contamination by non-monocyte cell types; $\mathbf{Z}_i^T\psi_u$ represent methylation chip and position random effects and are used to control for possible batch effects. The error term $\epsilon_i \sim MVN(0, \sigma^2 I_n)$ and is independent of the random effects. We obtained p -values for testing the null hypothesis $\psi_a = 0$ from the above model. We further applied the R/Bioconductor package BACON [8] to these p -values to further adjust for possible inflation using an empirical null distribution. Based on these marginal p -values, we selected top 2,000 CpG sites with the smallest p -values for our Bayesian multivariate analysis.

Besides the proposed methods, we also implement the other competing methods on the MESA data. HIMA identifies one CpG site in the gene region of *PCID2* as active mediator through the joint significance test with adjusted p -value = $6.3e-5$, and this single site has also been detected by PTG and GMM methods. We apply the Pathway Lasso and bi-Lasso on multiple permuted data, and notice that same active mediators with non-zero indirect effects have been picked out in both original and permuted data. Thus, the signals identified by Pathway Lasso and bi-Lasso are very probably false discoveries. For BAMA, the estimated PIPs over the 2,000 CpG sites are no more than 0.1, which does not provide strong evidence on the finding of active mediators.

	Full Sample (n, %)	Neighborhood Socioeconomic Disadvantage Mean (SD)	Body Mass Index (BMI) Mean (SD)
Full sample	1225 (100)	-0.32 (1.11)	29.5 (5.49)
Age			
55–65 years	462 (38)	-0.18 (0.96)	30.3 (6.02)
66–75 years	397 (32)	-0.30 (1.16)	30.1 (5.21)
76–85 years	300 (24)	-0.47 (1.15)	28.2 (4.65)
86–95 years	66 (5)	-0.67 (1.46)	26.6 (4.66)
Race/ethnic group			
Non-Hispanic white	580 (47)	-0.56 (1.18)	28.7 (5.40)
African-American	263 (22)	-0.16 (0.98)	30.5 (5.69)
Hispanic	382 (31)	-0.05 (1.00)	30.0 (5.32)
Gender			
Female	633 (52)	-0.24 (1.09)	30.1 (6.20)
Male	592 (48)	-0.40 (1.12)	28.9 (4.54)

Table S9: Characteristics of 1225 participants from MESA. %: proportion in the corresponding category. SD: standard deviation.

6 Detailed Description of LIFECODES Data

The LIFECODES prospective birth cohort enrolled approximately 1,600 pregnant women between 2006 and 2008 at the Brigham and Women’s Hospital in Boston, MA. Participants between 20 and 46 years of age were all at less than 15 weeks gestation at the initial study visit, and followed up to four visits (targeted at median 10, 18, 26, and 35 weeks gestation). At the initial study visit, questionnaires were administered to collect demographic and health-related information. Subjects’ urine and plasma samples were collected at each study visit. Among participants recruited in the LIFECODES cohort, 1,181 participants were

followed until delivery and had live singleton infants. The birth outcome, gestational age, was also recorded at delivery, and preterm birth was defined as delivery prior to 37 weeks gestation. This study received institutional review board (IRB) approval from the Brigham and Women’s Hospital and all participants provided written informed consent. All of the methods within this study were performed in accordance with the relevant guidelines and regulations approved by the IRB. Additional details regarding recruitment and study design can be found in [9, 10].

In this study, we focused on a subset of $n = 161$ individuals with their urine and plasma samples collected at one study visit occurring between 23.1 and 28.9 weeks gestation (median = 26 weeks). Characteristics of the subset sample is described in Table S10. Subjects’ urine samples were refrigerated (4°C) for a maximum of 2 hours before being processed and stored at -80°C . Approximately 10mL of blood was collected using ethylenediaminetetraacetic acid plasma tubes and temporarily stored at 4°C for less than 4 hours. Afterwards, blood was centrifuged for 20 minutes and stored at -80°C . Environmental exposure analytes were measured from urine samples by NSF International in Ann Arbor, MI, following the methods developed by the Centers for Disease Control (CDC) [11]. Those exposure analytes include phthalates, phenols and parabens, trace metals and polycyclic aromatic hydrocarbons. To adjust for urinary dilution, specific gravity (SG) levels were measured in each urine sample using a digital handheld refractometer (ATAGO Company Ltd., Tokyo, Japan), and was included as a covariate in regression models. Urine and plasma were subsequently analyzed for endogenous biomarkers, including 51 eicosanoids, five oxidative stress biomarkers and five immunological biomarkers in the present study. For a detailed description of the biomarkers that we analyzed and the media (urine or plasma) in which they were measured, please refer to [12].

	Full Sample (n = 161)	Preterm (<37 weeks gestation, n = 52)	Control (n = 109)
Age^a	32.7 (4.4)	32.1 (5.0)	33.0 (4.2)
BMI at Initial Visit^a	26.7 (6.4)	28.5 (7.6)	25.8 (5.6)
Race/ethnic group^b			
White	102 (63%)	29 (56%)	73 (67%)
African-American	18 (11%)	7 (13%)	11 (10%)
Other	41 (26%)	16 (31%)	25 (23%)
Gestational weeks^a	37.5 (3.1)	34.1 (3.2)	39.1 (1.1)

Table S10: Characteristics of all participants in the subset sample from the LIFECODES prospective birth cohort (n = 161). ^aContinuous variables presented as: mean (standard deviation). ^bCategorical variables presented as: count (proportion).

References

- [1] Donald B Rubin. “Randomization analysis of experimental data: The Fisher randomization test comment”. In: *Journal of the American Statistical Association* 75.371 (1980), pp. 591–593.
- [2] Donald B Rubin. “Comment: Which ifs have causal answers”. In: *Journal of the American Statistical Association* 81.396 (1986), pp. 961–962.
- [3] Tyler VanderWeele and Stijn Vansteelandt. “Mediation analysis with multiple mediators”. In: *Epidemiologic methods* 2.1 (2014), pp. 95–115.
- [4] Diane E Bild et al. “Multi-ethnic study of atherosclerosis: objectives and design”. In: *American journal of epidemiology* 156.9 (2002), pp. 871–881.
- [5] Yongmei Liu et al. “Methylomics of gene expression in human monocytes”. In: *Human molecular genetics* 22.24 (2013), pp. 5065–5074.
- [6] Pan Du, Warren A Kibbe, and Simon M Lin. “lumi: a pipeline for processing Illumina microarray”. In: *Bioinformatics* 24.13 (2008), pp. 1547–1548.
- [7] Yi-an Chen et al. “Discovery of cross-reactive probes and polymorphic CpGs in the Illumina Infinium HumanMethylation450 microarray”. In: *Epigenetics* 8.2 (2013), pp. 203–209.
- [8] Maarten van Iterson, Erik W van Zwet, and Bastiaan T Heijmans. “Controlling bias and inflation in epigenome-and transcriptome-wide association studies using the empirical null distribution”. In: *Genome biology* 18.1 (2017), p. 19.

- [9] Thomas F McElrath et al. “Longitudinal evaluation of predictive value for preeclampsia of circulating angiogenic factors through pregnancy”. In: *American journal of obstetrics and gynecology* 207.5 (2012), 407–e1.
- [10] Kelly K Ferguson et al. “Variability in urinary phthalate metabolite levels across pregnancy and sensitive windows of exposure for the risk of preterm birth”. In: *Environment international* 70 (2014), pp. 118–124.
- [11] Manori J Silva et al. “Quantification of 22 phthalate metabolites in human urine”. In: *Journal of Chromatography B* 860.1 (2007), pp. 106–112.
- [12] Max T Aung et al. “Prediction and associations of preterm birth and its subtypes with eicosanoid enzymatic pathways and inflammatory markers”. In: *Scientific reports* 9.1 (2019), pp. 1–17.