

# Compositional Mediation Model (CMM) for Binary Outcome

**Michael B. Sohn** & Hongzhe Li<sup>†</sup>

University of Rochester

Dept. Biostatistics and Computational Biology

February 7, 2019

<sup>†</sup>University of Pennsylvania

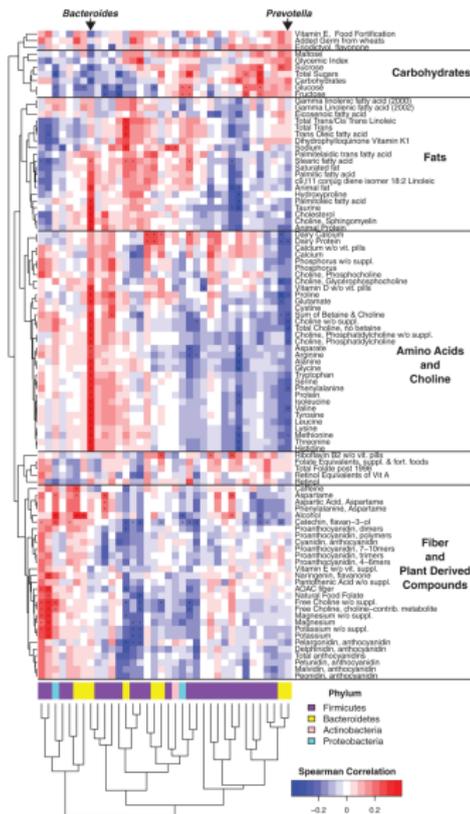
# COMBO Dataset (Wu, *et al.*, 2011 Science)

- Cross-sectional study Of diet and stool MicroBiOme composition (COMBO)

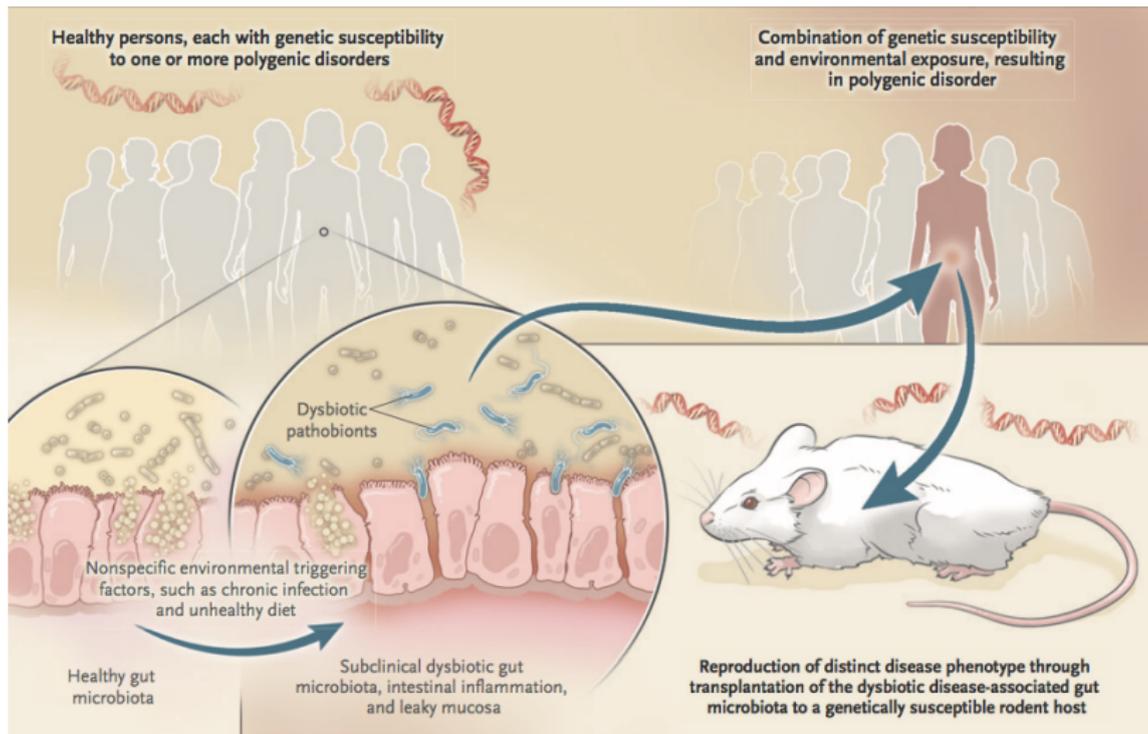
## • Data

- ▶ 98 healthy subjects/stool samples, not on antibiotics
- ▶ 16S rRNA gene sequences
- ▶ 87 genera appeared in at least one sample
- ▶ Nutrients (FFQ diet questionnaire) & demographic data such as BMI

## • Findings:

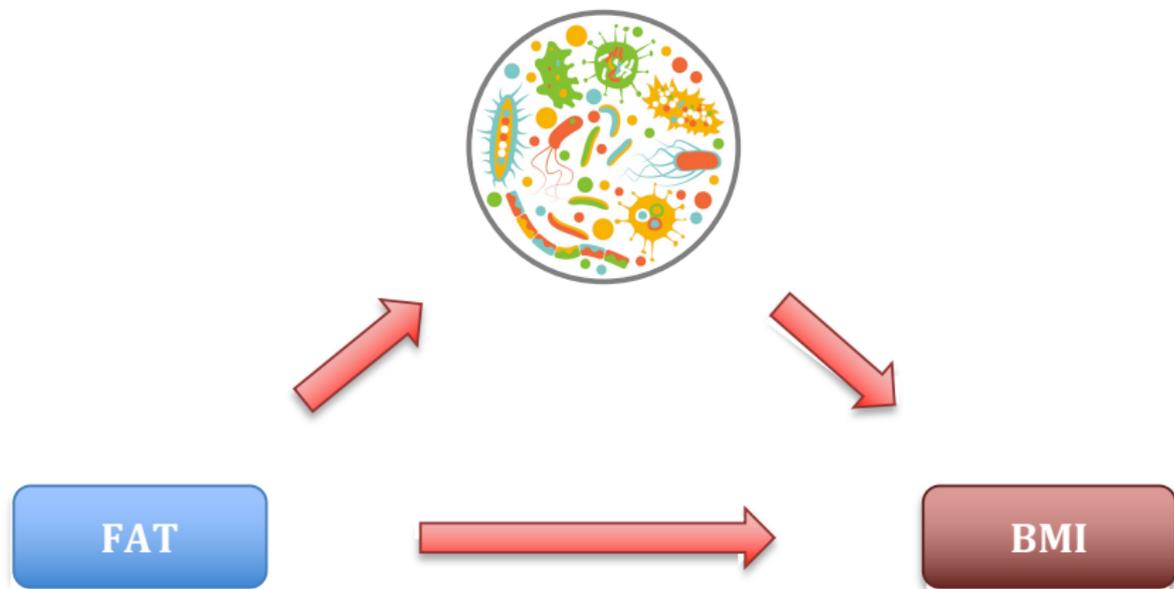


# Hypothesis of Pathogenesis Caused by Dysbiosis

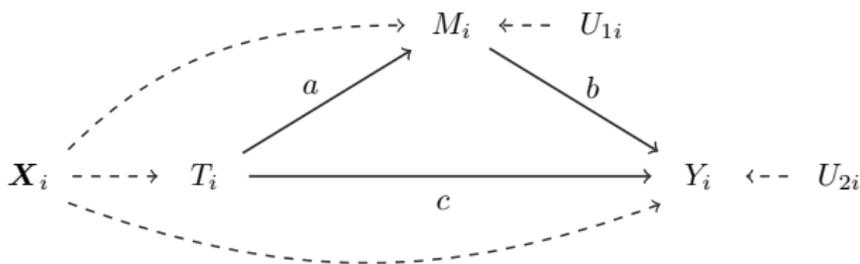


Source: N Engl J Med 2016;375:2369-79

## COMBO: Mediation Effect



# Mediation Analysis - Structural Equation Model (SEM)



$T_i$  - Treatment,  $M_i$  - Mediator,  $Y_i$  - Outcome,  $\mathbf{X}_i$  - Pretreatment Variables

$$M_i = a_0 + aT_i + \mathbf{h}^\top \mathbf{X}_i + U_{1i} \quad (1)$$

$$Y_i = c_0 + cT_i + bM_i + \mathbf{g}^\top \mathbf{X}_i + U_{2i} \quad (2)$$

By combining Eq. (1) and (2),

$$\begin{aligned} Y_i &= c_0 + cT_i + b(a_0 + aT_i + \mathbf{h}^\top \mathbf{X}_i + U_{1i}) + \mathbf{g}^\top \mathbf{X}_i + U_{2i} \\ &= c_0^* + (c + ab)T_i + \mathbf{g}^{*\top} \mathbf{X}_i + U_i^* \end{aligned}$$

# Mediation Analysis - Potential Outcomes Framework

Let  $T_i$  represent the binary treatment variable

**Causal Direct Effect:**  $\zeta(t) = \mathbb{E}[Y_i(1, M_i(t)|\mathbf{X}_i) - Y_i(0, M_i(t)|\mathbf{X}_i)]$

**Causal Indirect Effect:**  $\delta(t) = \mathbb{E}[Y_i(t, M_i(1)|\mathbf{X}_i) - Y_i(t, M_i(0)|\mathbf{X}_i)]$

Necessary Assumptions:

- Stable Unit Treatment Value Assumption (SUTVA)
- Sequential Ignorability Assumption,

$$\{Y_i(t', m), M_i(t)\} \perp\!\!\!\perp T_i | \mathbf{X}_i = \mathbf{x},$$

$$Y_i(t', m) \perp\!\!\!\perp M_i(t) | T_i = t, \mathbf{X}_i = \mathbf{x},$$

where  $0 < Pr(T_i = t | \mathbf{X}_i = \mathbf{x})$  and  $0 < Pr(M_i(t) = m | T_i = t, \mathbf{X}_i = \mathbf{x})$  for  $t = 0, 1$ .

With the necessary assumptions,

$$\zeta(t) = \int \mathbb{E}(Y_i | M_i, T_i, \mathbf{X}_i) [dF_{M_i | T_i=1, \mathbf{X}_i}(m) - dF_{M_i | T_i=0, \mathbf{X}_i}(m)] dF_{\mathbf{X}_i}(\mathbf{x})$$

$$\delta(t) = \int [\mathbb{E}(Y_i | M_i, T_i = 1, \mathbf{X}_i) - \mathbb{E}(Y_i | M_i, T_i = 0, \mathbf{X}_i)] dF_{M_i | T_i, \mathbf{X}_i}(m) dF_{\mathbf{X}_i}(\mathbf{x})$$

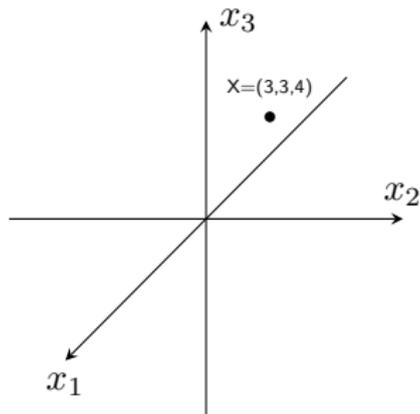
# Compositional Data Analysis

Compositional data:

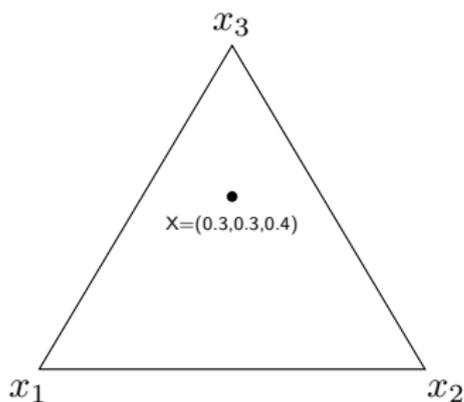
- Relative information
- Proportions or percentages of a whole

Unit-sum constraint: sum of proportions = 1

Euclidean Space



Simplex Space



## Subcompositional Coherence

**Principle of subcompositional coherence:** analysis concerning a subset of components must not depend on excluded components

**Example:** Scientists A and B record the composition of soil samples:

A records **animal**, **vegetable**, **mineral**, and **water**.

B records **animal**, **vegetable**, and **mineral** after drying the sample.

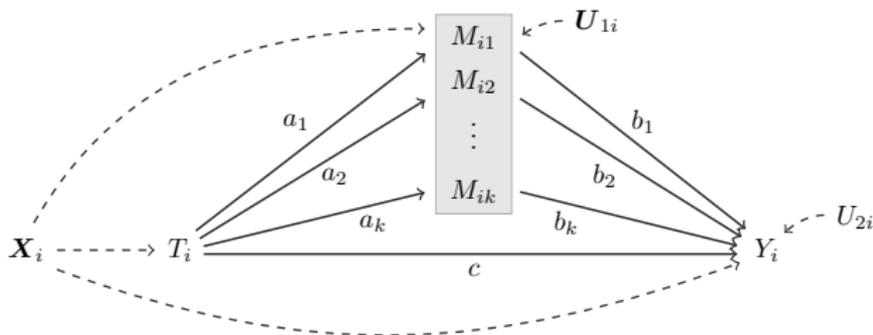
Both are absolutely accurate. [adapted from Aitchison, 2005]

Sample A	$x_1$	$x_2$	$x_3$	$x_4$
1	0.1	0.2	0.1	0.6
2	0.2	0.1	0.2	0.5
3	0.3	0.3	0.1	0.3

Sample B	$x_1$	$x_2$	$x_3$
1	0.25	0.50	0.25
2	0.40	0.20	0.40
3	0.43	0.43	0.14

Corr A	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	1.00	0.50	0.00	-0.98
$x_2$		1.00	-0.87	-0.65
$x_3$			1.00	0.19

Corr B	$x_1$	$x_2$	$x_3$
$x_1$	1.00	-0.57	-0.05
$x_2$		1.00	-0.79
$x_3$			1.00



Compositional operators (Aitchison, 1986; Billheimer, *et al.* 2001):

$$\mathbf{m} \oplus \mathbf{a} = \left( \frac{m_1 a_1}{\sum_{j=1}^k m_k a_k}, \dots, \frac{m_k a_k}{\sum_{j=1}^k m_k a_k} \right)^\top; \quad \mathbf{m}^z = \left( \frac{m_1^z}{\sum_{j=1}^k m_k^z}, \dots, \frac{m_k^z}{\sum_{j=1}^k m_k^z} \right)^\top$$

Compositional mediation model:

$$\mathbf{M}_i = \left( \mathbf{m}_0 \oplus \mathbf{a}^{T_i} \bigoplus_{r=1}^{n_x} \mathbf{h}_r^{X_{ri}} \right) \oplus \mathbf{U}_{1i}$$

$$Y_i = c_0 + cT_i + \mathbf{b}^\top (\log \mathbf{M}_i) + \mathbf{g}^\top \mathbf{X}_i + U_{2i}, \quad \text{subject to } \mathbf{1}_k^\top \mathbf{b} = 0$$

Necessary assumptions:

- SUTVA
- Sequential Ignorability Assumption

Under the potential outcomes framework

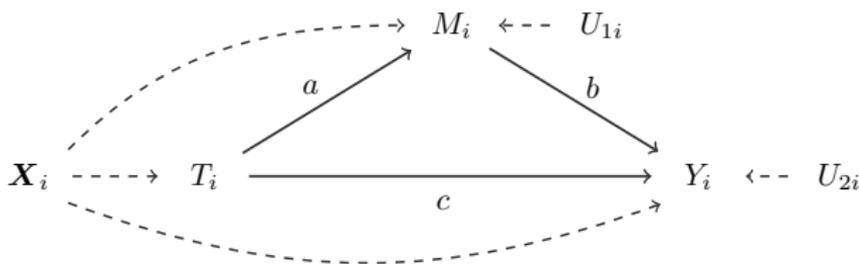
- Expected Causal Direct Effect

$$\begin{aligned}\zeta(t) &= \mathbb{E}[Y_i(1, \log \mathbf{M}_i(t) | \mathbf{X}_i(t)) - Y_i(0, \log \mathbf{M}_i(t) | \mathbf{X}_i(t))] \\ &= \mathbf{c}\end{aligned}$$

- Expected Causal Indirect Effect

$$\begin{aligned}\delta(t) &= \mathbb{E}[Y_i(t, \log \mathbf{M}_i(1) | \mathbf{X}_i(t)) - Y_i(t, \log \mathbf{M}_i(0) | \mathbf{X}_i(t))] \\ &= (\log \mathbf{a})^\top \mathbf{b}\end{aligned}$$

## Binary Outcome Under SEM



$T_i$  - Treatment,  $M_i$  - Mediator,  $Y_i$  - Binary outcome,  $\mathbf{X}_i$  - Pretreatment Variables

$$M_i = a_0 + aT_i + \mathbf{h}^\top \mathbf{X}_i + U_{1i} \quad (3)$$

$$Y_i = \mathbf{1}\{Y_i^* > 0\}, \text{ where } Y_i^* = c_0 + cT_i + bM_i + \mathbf{g}^\top \mathbf{X}_i + U_{2i} \quad (4)$$

By combining Eq. (3) and (4),

$$\begin{aligned} Y_i^* &= c_0 + cT_i + b(a_0 + aT_i + \mathbf{h}^\top \mathbf{X}_i + U_{1i}) + \mathbf{g}^\top \mathbf{X}_i + U_{2i} \\ &= c_0^* + (\mathbf{c} + \mathbf{ab})T_i + \mathbf{g}^{*\top} \mathbf{X}_i + U_i^* \end{aligned}$$

Assumptions: SUTVA, Sequential ignorability

Estimation of Causal Direct  $\zeta$  and Indirect Effects  $\delta$ :

- Logit Model
  - Complex numerical integration required
  - Odds ratios with rare outcomes:  $\log OR_\zeta \approx c$ ;  $\log OR_\delta \approx ab$
- Probit Model

$$\zeta = \mathbb{E} \left\{ \Phi \left( \frac{c + f_1}{\sqrt{\sigma^2 b^2 + 1}} \right) - \Phi \left( \frac{f_1}{\sqrt{\sigma^2 b^2 + 1}} \right) \right\}$$

$$\delta = \mathbb{E} \left\{ \Phi \left( \frac{ab + f_2}{\sqrt{\sigma^2 b^2 + 1}} \right) - \Phi \left( \frac{f_2}{\sqrt{\sigma^2 b^2 + 1}} \right) \right\}$$

where  $f_1 = c_0 + a_0 b + b(\mathbf{h} + \mathbf{g})^\top \mathbf{X}_i$  and  $f_2 = c_0 + c + a_0 b + b(\mathbf{h} + \mathbf{g})^\top \mathbf{X}_i$

## CMM for Binary Outcome (Probit Model)

Compositional mediation model for the binary outcome:

$$\mathbf{M}_i = \left( \mathbf{m}_0 \oplus \mathbf{a}^{T_i} \bigoplus_{r=1}^{n_x} \mathbf{h}_r^{X_{ri}} \right) \oplus \mathbf{U}_{1i}$$

$$Y_i = 1\{c_0 + cT_i + \mathbf{b}^\top (\log \mathbf{M}_i) + \mathbf{g}^\top \mathbf{X}_i + U_{2i} > 0\}, \text{ subject to } \mathbf{1}_k^\top \mathbf{b} = 0$$

where  $\mathbf{U}_{1i} \sim LN(0, \Sigma)$  and  $U_{2i} \sim N(0, 1)$ .

### Expected Causal Direct and Indirect Effects:

$$\zeta(\tau) = \mathbb{E} \left\{ \Phi \left( \frac{ct + f_\zeta(\tau, \mathbf{X}_i)}{\sqrt{\mathbf{b}_{-k}^\top \Sigma \mathbf{b}_{-k} + 1}} \right) - \Phi \left( \frac{ct' + f_\zeta(\tau, \mathbf{X}_i)}{\sqrt{\mathbf{b}_{-k}^\top \Sigma \mathbf{b}_{-k} + 1}} \right) \right\}$$

$$\delta(\tau) = \mathbb{E} \left\{ \Phi \left( \frac{(\log \mathbf{a})^\top \mathbf{b}t + f_\delta(\tau, \mathbf{X}_i)}{\sqrt{\mathbf{b}_{-k}^\top \Sigma \mathbf{b}_{-k} + 1}} \right) - \Phi \left( \frac{(\log \mathbf{a})^\top \mathbf{b}t' + f_\delta(\tau, \mathbf{X}_i)}{\sqrt{\mathbf{b}_{-k}^\top \Sigma \mathbf{b}_{-k} + 1}} \right) \right\}$$

where  $f_\zeta(\tau, \mathbf{x}) = c_0 + \mathbf{b}^\top (\log \mathbf{m}_0 + \tau \log \mathbf{a} + \sum_{r=1}^{n_x} x_r \log \mathbf{h}_r) + \mathbf{g}^\top \mathbf{x}$ ;

$f_\delta(\tau, \mathbf{x}) = c_0 + c\tau + \mathbf{b}^\top (\log \mathbf{m}_0 + \sum_{r=1}^{n_x} x_r \log \mathbf{h}_r) + \mathbf{g}^\top \mathbf{x}$ .

Optimization problem in a simplex space:

$$(\widehat{\mathbf{m}}_0, \widehat{\mathbf{a}}, \widehat{\mathbf{h}}_1, \dots, \widehat{\mathbf{h}}_{n_x}) = \underset{\mathbf{m}_0, \mathbf{a}, \mathbf{h}_r \in \mathbb{S}^{k-1}}{\operatorname{argmin}} \sum_{i=1}^n \left\| \mathbf{M}_i \ominus \left( \mathbf{m}_0 \oplus \mathbf{a}^{T_i} \bigoplus_{r=1}^{n_x} \mathbf{h}_r^{X_{r,i}} \right) \right\|^2$$

where

$$\mathbf{m} \ominus \mathbf{a} = \left( \frac{m_1 a_1^{-1}}{\sum_{j=1}^k m_j a_j^{-1}}, \frac{m_2 a_2^{-1}}{\sum_{j=1}^k m_j a_j^{-1}}, \dots, \frac{m_k a_k^{-1}}{\sum_{j=1}^k m_j a_j^{-1}} \right)$$

$$\|\mathbf{m}\| = \langle \mathbf{m}, \mathbf{m} \rangle^{1/2} = \operatorname{alr}(\mathbf{m})^\top \mathcal{N}^{-1} \operatorname{alr}(\mathbf{m})$$

$$\operatorname{alr}(\mathbf{m}) = \left( \log \frac{m_1}{m_k}, \log \frac{m_2}{m_k}, \dots, \log \frac{m_{k-1}}{m_k} \right)^\top$$

$$\mathcal{N}^{-1} = \mathcal{I}_{k-1} - \frac{1}{k} \mathbf{1}_{k-1} \mathbf{1}_{k-1}^\top$$

# Estimator for Parameters in Probit Regression

Let  $\eta_i = 2y_i - 1$ ,  $\mathbf{z}_i = (1, t_i, \log(\mathbf{m}_i)^\top, \mathbf{x}_i^\top)^\top$ , and  $\boldsymbol{\beta} = (c_0, c, \mathbf{b}^\top, \mathbf{g}^\top)^\top$

$$\widehat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^n \log \Phi(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad \text{s.t. } \mathbf{1}_k^\top \mathbf{b} = 0 \quad (5)$$

**Alternative optimization problem:**

$$\widehat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|\Xi^{1/2}(\mathbf{u} - Z\boldsymbol{\beta})\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad \text{s.t. } \mathbf{1}_k^\top \mathbf{b} = 0, \quad (6)$$

where  $\Xi$  is the  $n \times n$  diagonal matrix with  $\Xi_{ii} = \xi_i(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta}^*) [ \mathbf{z}_i^\top \boldsymbol{\beta}^* + \xi_i(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta}^*) ]$ ,  $\xi_i(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta}) = \eta_i \phi(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta}) / \Phi(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta})$ ,  $\mathbf{u} = Z\boldsymbol{\beta}_0 + (\Xi)^{-1} \boldsymbol{\xi}$ ,  $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n)^\top$ , and  $\boldsymbol{\xi} = (\xi_1(\eta_1 \mathbf{z}_1^\top \boldsymbol{\beta}_0), \dots, \xi_n(\eta_n \mathbf{z}_n^\top \boldsymbol{\beta}_0))^\top$

## Proposed method: IRLS-CDMM

$$\begin{aligned}\widehat{\boldsymbol{\beta}}^{(\ell)} &= \operatorname{argmin}_{\boldsymbol{\beta}} \left\{ \frac{1}{2n} \left\| \Xi^{(\ell-1)1/2} (\mathbf{u}^{(\ell-1)} - Z\boldsymbol{\beta}) \right\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad \text{s.t. } \mathbf{1}_k^\top \boldsymbol{\beta}^{(\ell)} = 0, \\ &= \operatorname{argmin}_{\boldsymbol{\beta}} \left\{ \frac{1}{2n} \left\| \Xi^{(\ell-1)1/2} (\mathbf{u}^{(\ell-1)} - \widetilde{Z}\boldsymbol{\beta}) \right\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad \text{s.t. } \boldsymbol{\iota}^\top \boldsymbol{\beta}^{(\ell)} = 0,\end{aligned}$$

where  $\widetilde{Z} = Z(I_p - \boldsymbol{\iota}\boldsymbol{\iota}^\top/k)$  and  $\boldsymbol{\iota} = (0, 0, 1, \dots, 1, 0, \dots, 0)^\top$ .

## Algorithm: IRLS-CDMM with Augmented Lagrangian Method

*Step 1.* Initialize  $\boldsymbol{\beta}^{(0)}$ ,  $\alpha^{(0)}$

*Step 2.* Update  $\beta_j^{(\ell+1)}$  until convergence

*Step 3.* Update  $\Xi^{(\ell+1)}$  and  $\mathbf{u}^{(\ell+1)}$  by minimizing  $\sum_{i=1}^n q(\eta_i \mathbf{z}_i^\top \boldsymbol{\beta})$

*Step 4.* Update  $\alpha^{(k+1)}$

Debiased Estimator of (6):

$$\widehat{\beta}_{db} = \widehat{\beta} + \frac{1}{n} \widetilde{M} \widetilde{Z}^\top \Xi(\mathbf{u} - \widetilde{Z} \widehat{\beta}),$$

where  $\widetilde{M} = (I_p - \boldsymbol{\mu}^\top/k)M$  and  $M = (\mathbf{m}_1, \dots, \mathbf{m}_p)^\top$  is a solution of the convex problem (Javanmard and Montanari, 2014; Shi, et. al., 2016):

$$\min \mathbf{m}_j^\top \widehat{\Sigma} \mathbf{m}_j \quad \text{s.t.} \quad \|\widehat{\Sigma} \mathbf{m}_j - (I_p - \boldsymbol{\mu}^\top/k) \mathbf{e}_j\|_\infty \leq \gamma, \quad j = 1, \dots, p,$$

where  $\mathbf{e}_j$  is the  $j^{\text{th}}$  natural basis and  $\gamma$  is some constant.

**Theorem:** For an  $s$ -sparse  $\beta$ , under some regularity conditions,

$$\sqrt{n}(\widehat{\beta}_{db} - \beta) = R + \Delta, \quad \mathbb{E}(R|Z) = \mathbf{0}, \quad \|\Delta\|_\infty \rightarrow 0$$

Null Hypothesis for the expected causal direct and indirect effects:

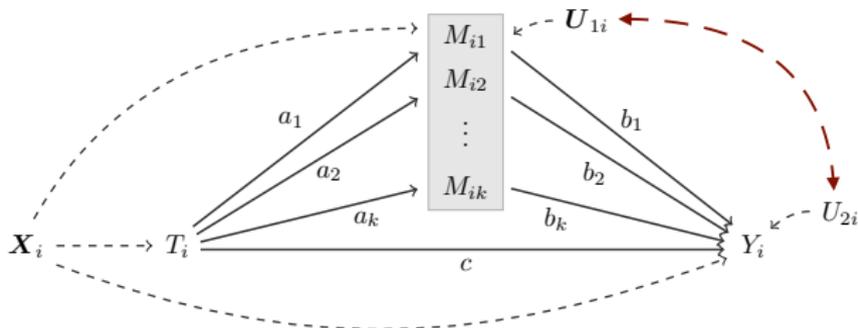
$$H_0 : \zeta(\tau) = 0 \quad \text{vs.} \quad H_1 : \zeta(\tau) \neq 0.$$

$$H_0 : \delta(\tau) = 0 \quad \text{vs.} \quad H_1 : \delta(\tau) \neq 0.$$

Testing Procedure (Non-parametric Bootstrap):

1. Randomly select  $n$  samples from the original  $n$  samples with replacement
2. Estimate  $\zeta_b(\tau)$  and  $\delta_b(\tau)$
3. Repeat Steps 1 and 2 to construct sampling distributions of  $\zeta_b(\tau)$  and  $\delta_b(\tau)$
4. Construct percentile bootstrap confidence intervals for  $\zeta_b(\tau)$  and  $\delta_b(\tau)$

# Sensitivity Analysis for Binary Outcome



Probit regression:  $Y_i = 1\{\tilde{c}_0 + \tilde{c}T_i + \tilde{\mathbf{g}}^\top \mathbf{X}_i + U_{0i} > 0\}$ , where  $U_{0i} \sim N(0, 1)$

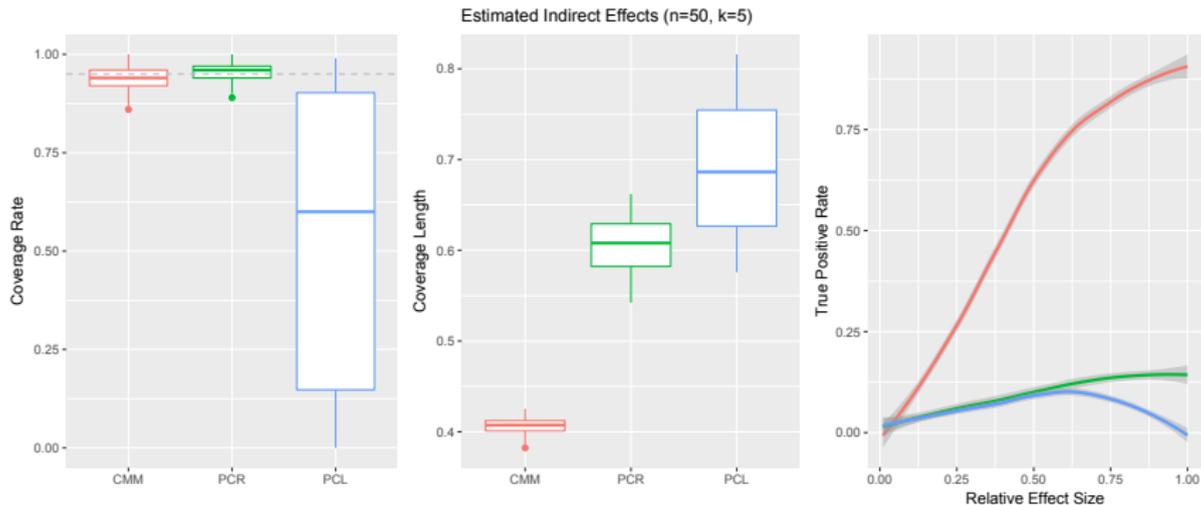
Expected causal indirect effect given  $\boldsymbol{\rho} = \text{corr}(\text{alt}(U_{1i}), U_{2i})$ :

$$\delta_\rho(\tau) = \mathbb{E} \left\{ \Phi \left( \tilde{f}_\delta(\tau) + \frac{(\log \mathbf{a})^\top \mathbf{b}_\rho (t - \tau)}{\Psi(\boldsymbol{\rho}, \mathbf{b}_\rho, \Sigma)} \right) - \Phi \left( \tilde{f}_\delta(\tau) + \frac{(\log \mathbf{a})^\top \mathbf{b}_\rho (t' - \tau)}{\Psi(\boldsymbol{\rho}, \mathbf{b}_\rho, \Sigma)} \right) \right\},$$

where  $\tilde{f}_\delta(\tau) = \tilde{c}_0 + \tilde{c}\tau + \tilde{\mathbf{g}}^\top \mathbf{x}_i$ ,  $\Psi(\boldsymbol{\rho}, \mathbf{b}_\rho, \Sigma) = [(\mathbf{b}_\rho)_{-k}^\top \Sigma (\mathbf{b}_\rho)_{-k} + 2\boldsymbol{\rho}^\top D (\mathbf{b}_\rho)_{-k} + 1]^{1/2}$ , and  $D$  is a diagonal matrix with  $\text{diag}(\Sigma^{1/2})$ .

# Performance Evaluation I ( $\alpha = 0.05$ )

Binary treatment ( $t = 1, t' = 0$ );  $\mathbf{a} = (20, 10, 5, 2, \mathbf{1}_{k-4}^\top)^\top / (20, 10, 5, 2, \mathbf{1}_{k-4}^\top)^\top \mathbf{1}_k$ ;  
 $\mathbf{b} = (0.5, -0.5, 0.5, -0.5, \mathbf{0}_{k-4}^\top)^\top$ ;  $(\log \mathbf{a})^\top \mathbf{b} \times \text{Effect Size}$

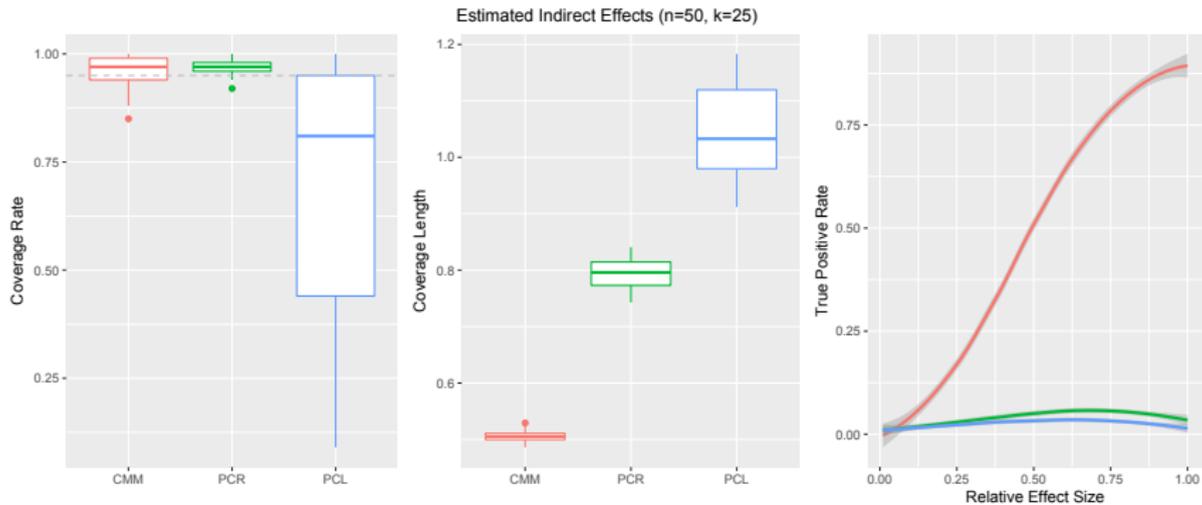


CMM: Proposed compositional mediation model

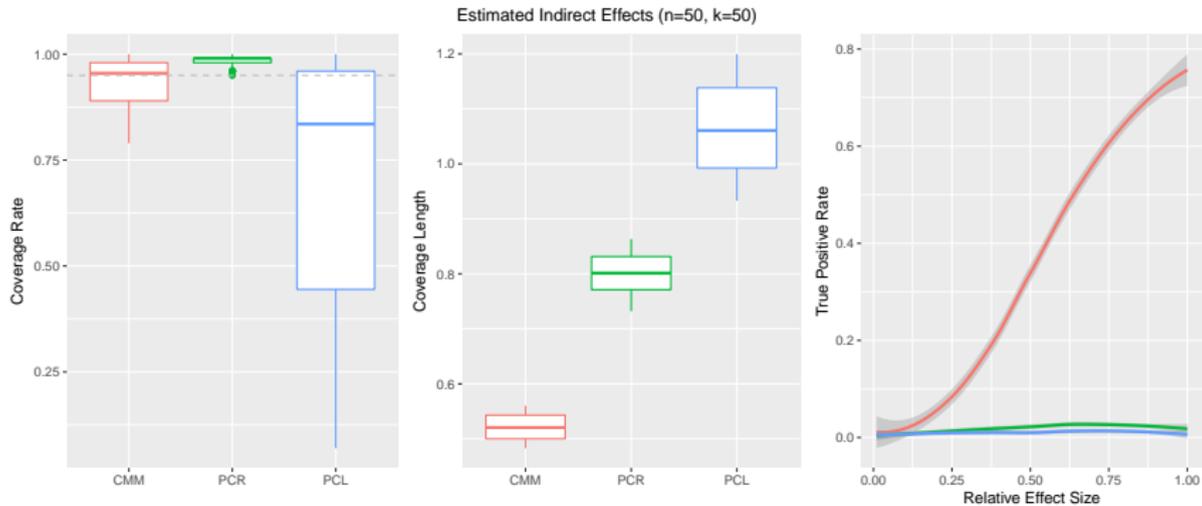
PCR: Principal components of compositional variables under POF

PCL: Principal components of compositional variables under SEM

# Performance Evaluation II ( $\alpha = 0.05$ )



# Performance Evaluation III ( $\alpha = 0.05$ )

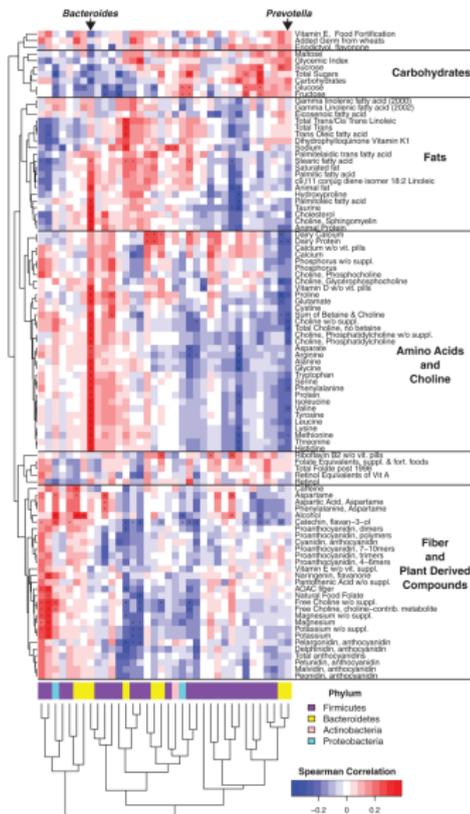
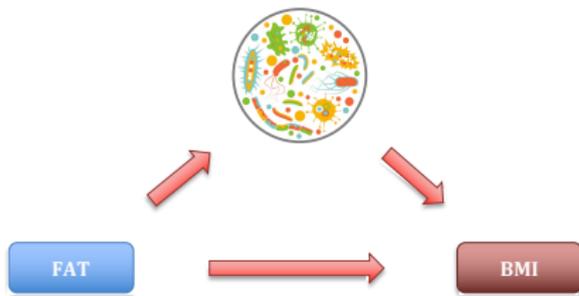


# COMBO Dataset

## • Data

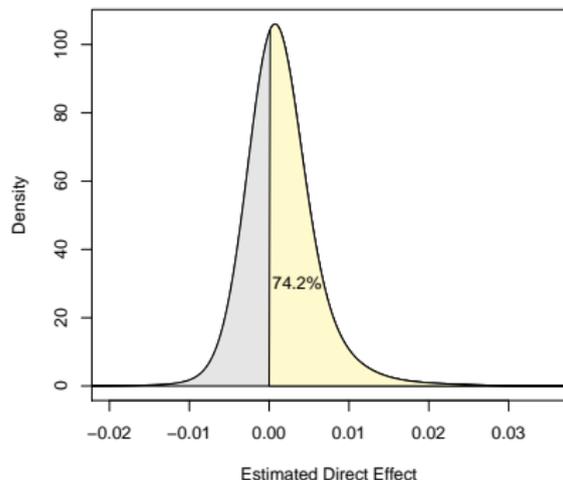
- ▶ 98 healthy subjects
- ▶ Fat intake as treatment
- ▶ 45 genera as compositional mediators
- ▶ Dichotomized BMI at 25

## • Interest:

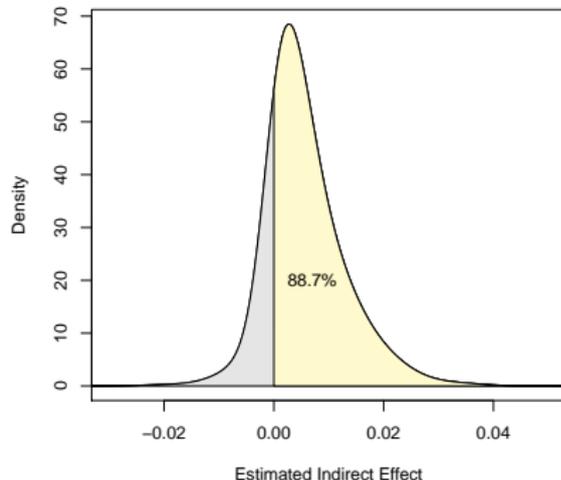


# Fat intake, Microbiome, and Obesity (COMBO)

Bootstrap Distribution of DE



Bootstrap Distribution of IDE



DE	TIDE
0.002 (-0.003, 0.011)	0.008 (-0.005, 0.023)
20%	80%

⇒ About 80% of effect of fat intake on dichotomized BMI is mediated through gut microbiome

Many thanks to

- Hongzhe Li, PhD