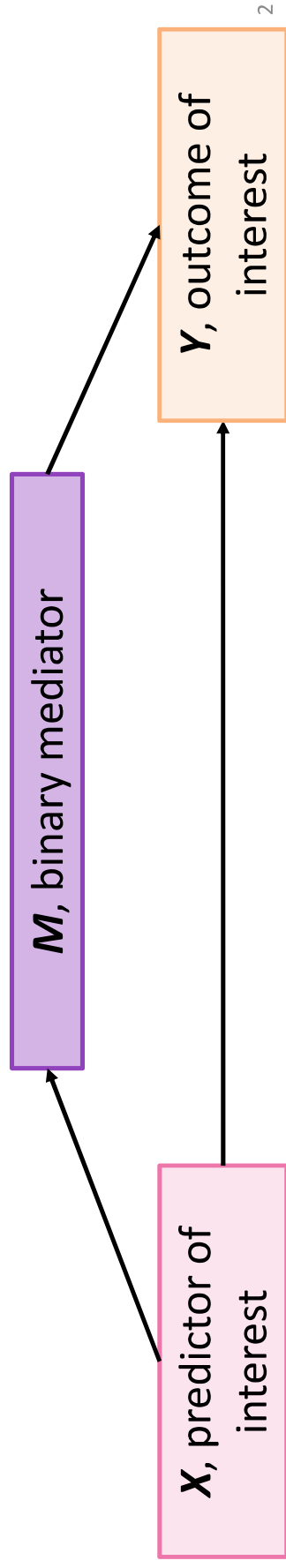# Effect estimation in the presence of a misclassified binary mediator

Kimberly A. H. Webb and Martin T. Wells

Cornell University
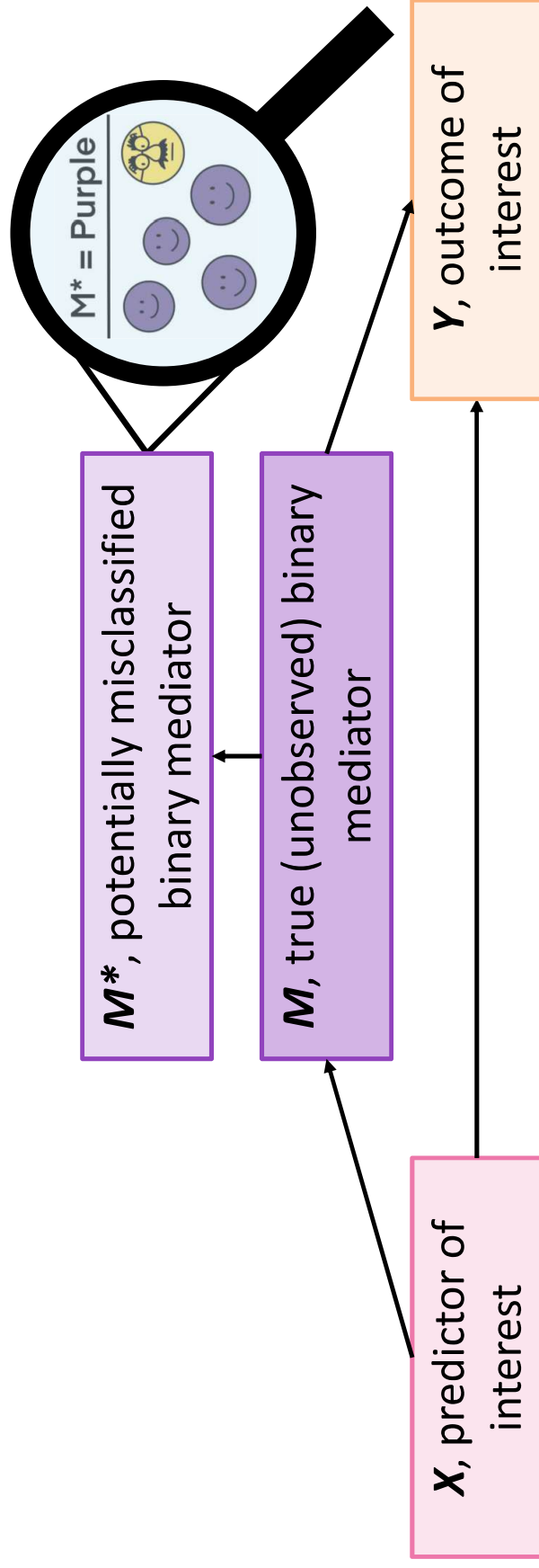
# Problem setting

- **Mediation analysis** quantifies the effect of an **exposure (X)** on an **outcome (Y)**, mediated by some **intermediate (M)**.



*M*, binary mediator

*Y*, outcome of interest

*X*, predictor of interest

# Problem setting

- **Mediation analysis** quantifies the effect of an <span style="color:#e8659e">exposure (X)</span> on an <span style="color:#f5a623">outcome (Y)</span>, mediated by some <span style="color:#8e44ad">intermediate (M)</span>.

- Measure <span style="color:#8e44ad">M</span> using an instrument that is not always accurate, and obtain <span style="color:#8e44ad">M*</span>.

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**Y**, outcome of interest
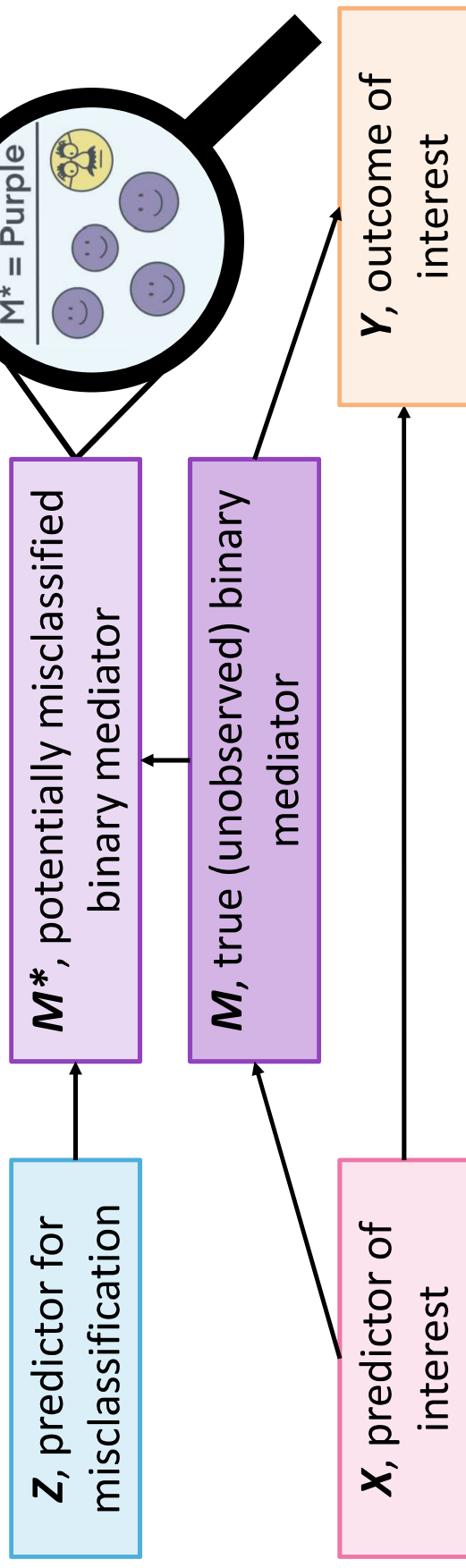
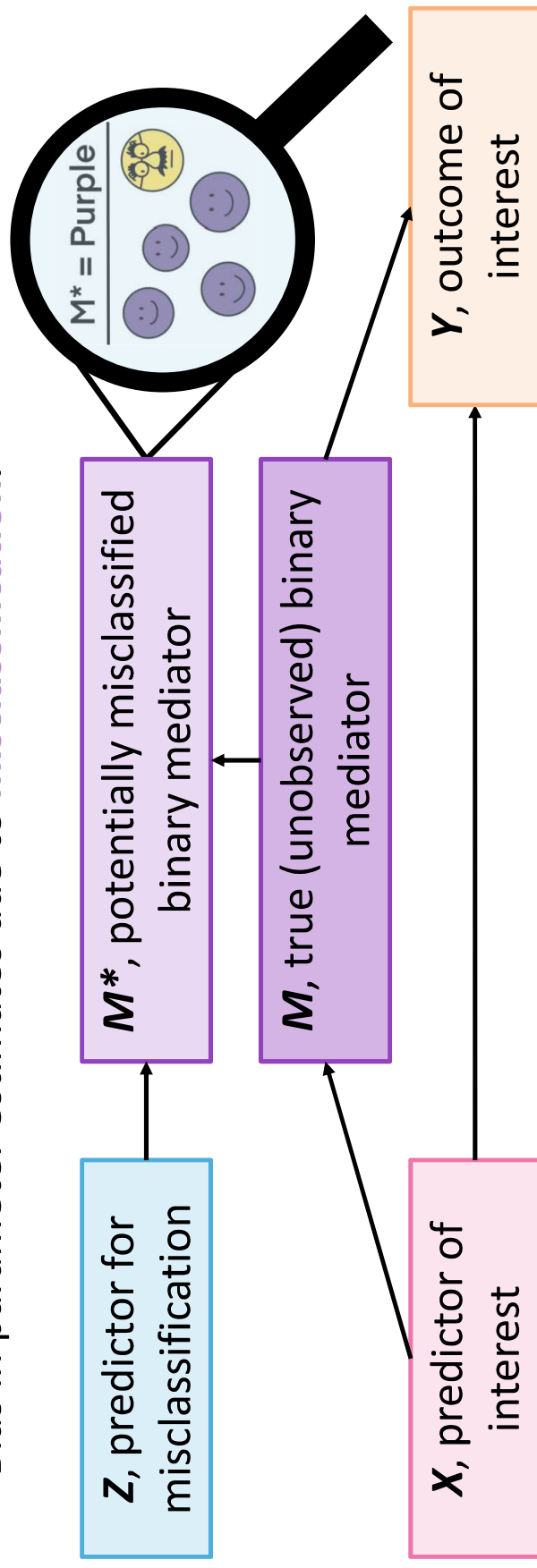**X**, predictor of interest

M* = Purple

# Problem setting

- **Mediation analysis** quantifies the effect of an **exposure (X)** on an **outcome (Y)**, mediated by some **intermediate (M)**.

- Measure **M** using an instrument that is not always accurate, and obtain **M\***.

- **Z** is related to the **misclassification mechanism**.



$M^*$, potentially misclassified binary mediator

$M$, true (unobserved) binary mediator

$Y$, outcome of interest

$Z$, predictor for misclassification

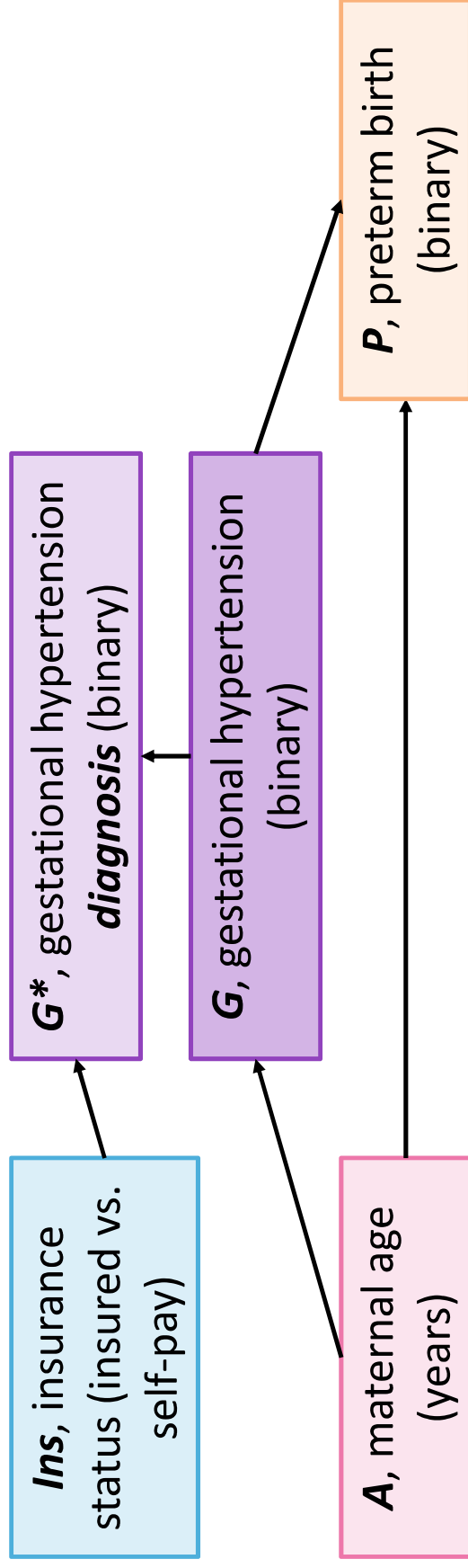$X$, predictor of interest

$M^* = $ Purple

# Problem setting

- **Challenges:**
  - Misclassification is covariate-dependent.
  - No gold standard labels.
  - Bias in parameter estimates due to misclassification.



**Z**, predictor for misclassification

**M\***, potentially misclassified binary mediator

M* = Purple

**M**, true (unobserved) binary mediator

**X**, predictor of interest

**Y**, outcome of interest

# Problem setting

- **Example:**
  - Does **gestational hypertension** mediate the association between **maternal age** and **preterm birth**, after accounting for potential **misdiagnosis of gestational hypertension** based on **patient insurance status**?
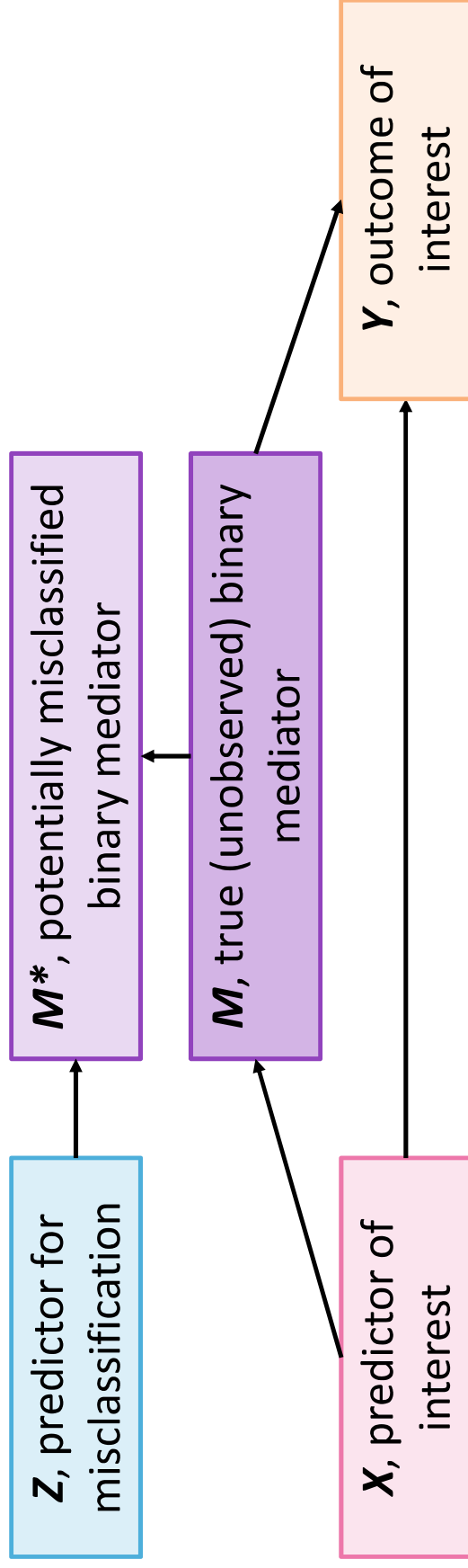
**Ins**, insurance status (insured vs. self-pay)

**G\***, gestational hypertension **diagnosis** (binary)

**G**, gestational hypertension (binary)

**A**, maternal age (years)

**P**, preterm birth (binary)

# Model

**True mediator model:**

**Observed mediator model:**

**Outcome model:**

**Z**, predictor for misclassification

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**X**, predictor of interest

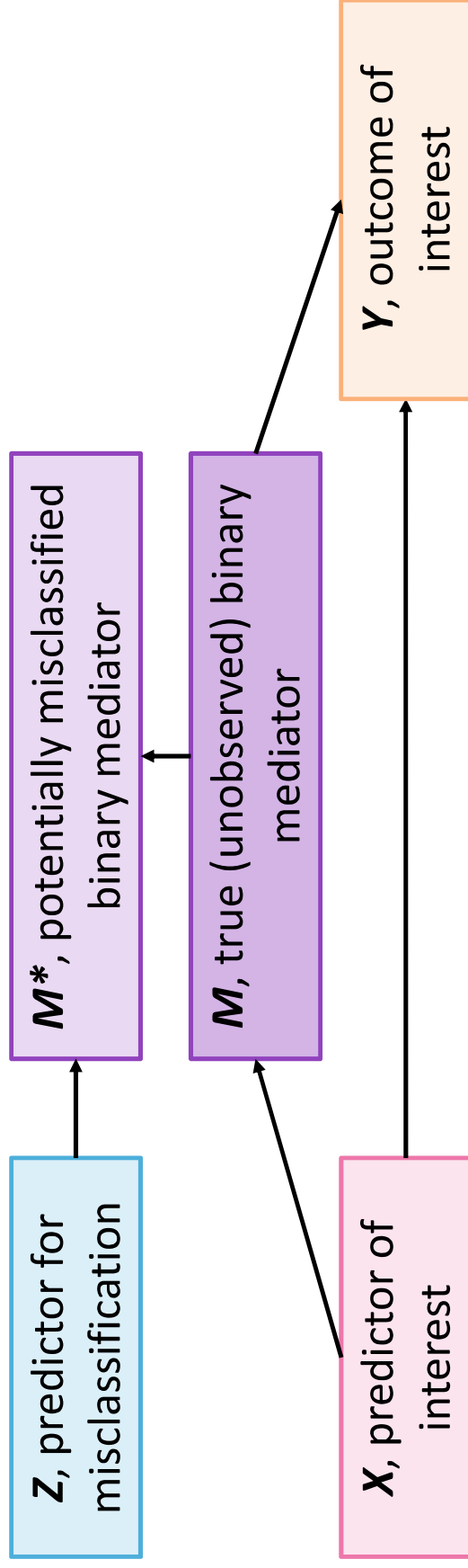**Y**, outcome of interest

# Model

**True mediator model:** $\text{logit}\{P(M = 1 | X, C; \boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:**

**Outcome model:**



**Z**, predictor for misclassification

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**X**, predictor of interest

**Y**, outcome of interest

# Model

**True mediator model:** $\text{logit}\{P(M = 1 | X, C; \boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_C C$

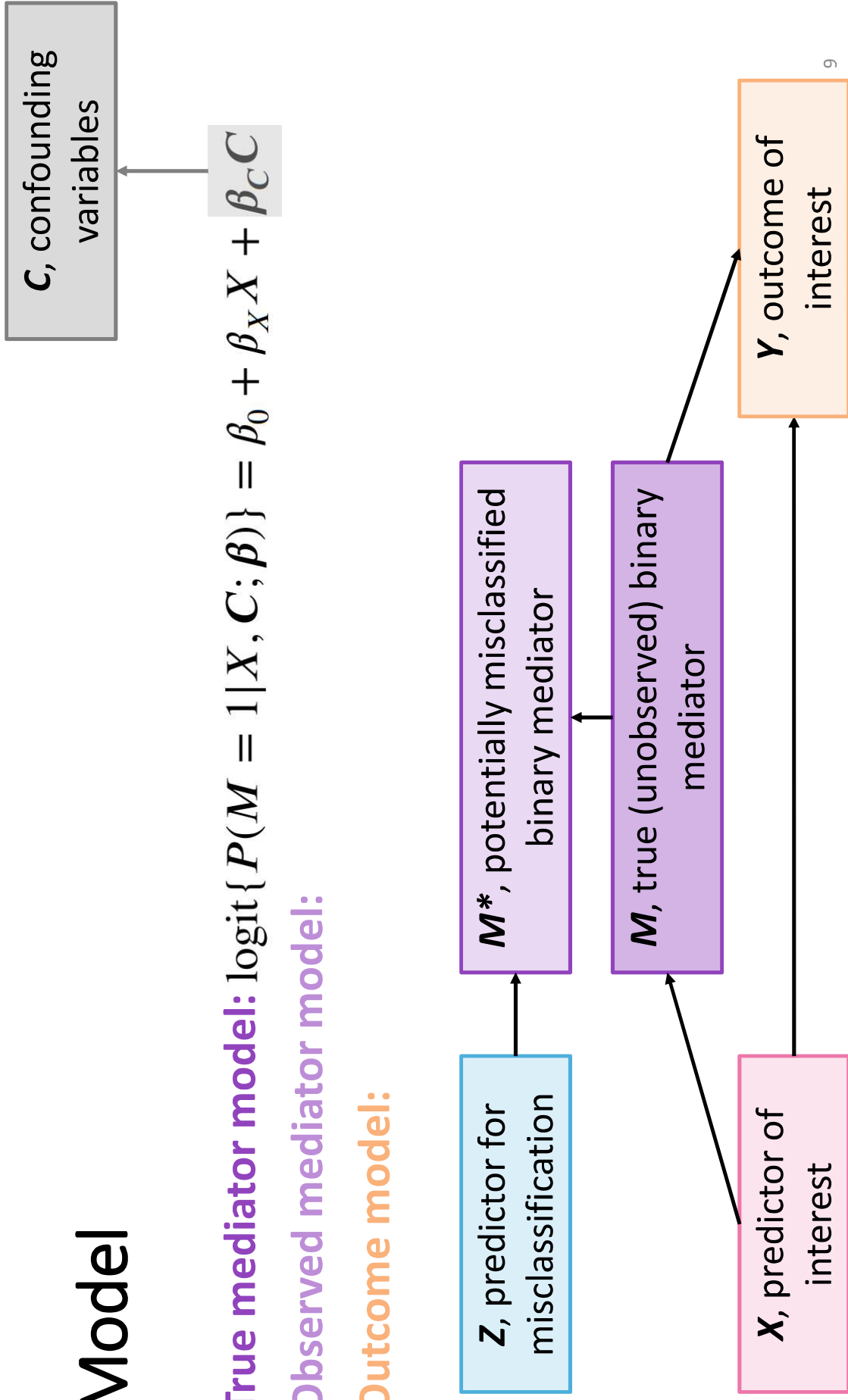**Observed mediator model:**

**Outcome model:**

$C$, confounding variables

$M^*$, potentially misclassified binary mediator

$M$, true (unobserved) binary mediator

$Y$, outcome of interest

$Z$, predictor for misclassification
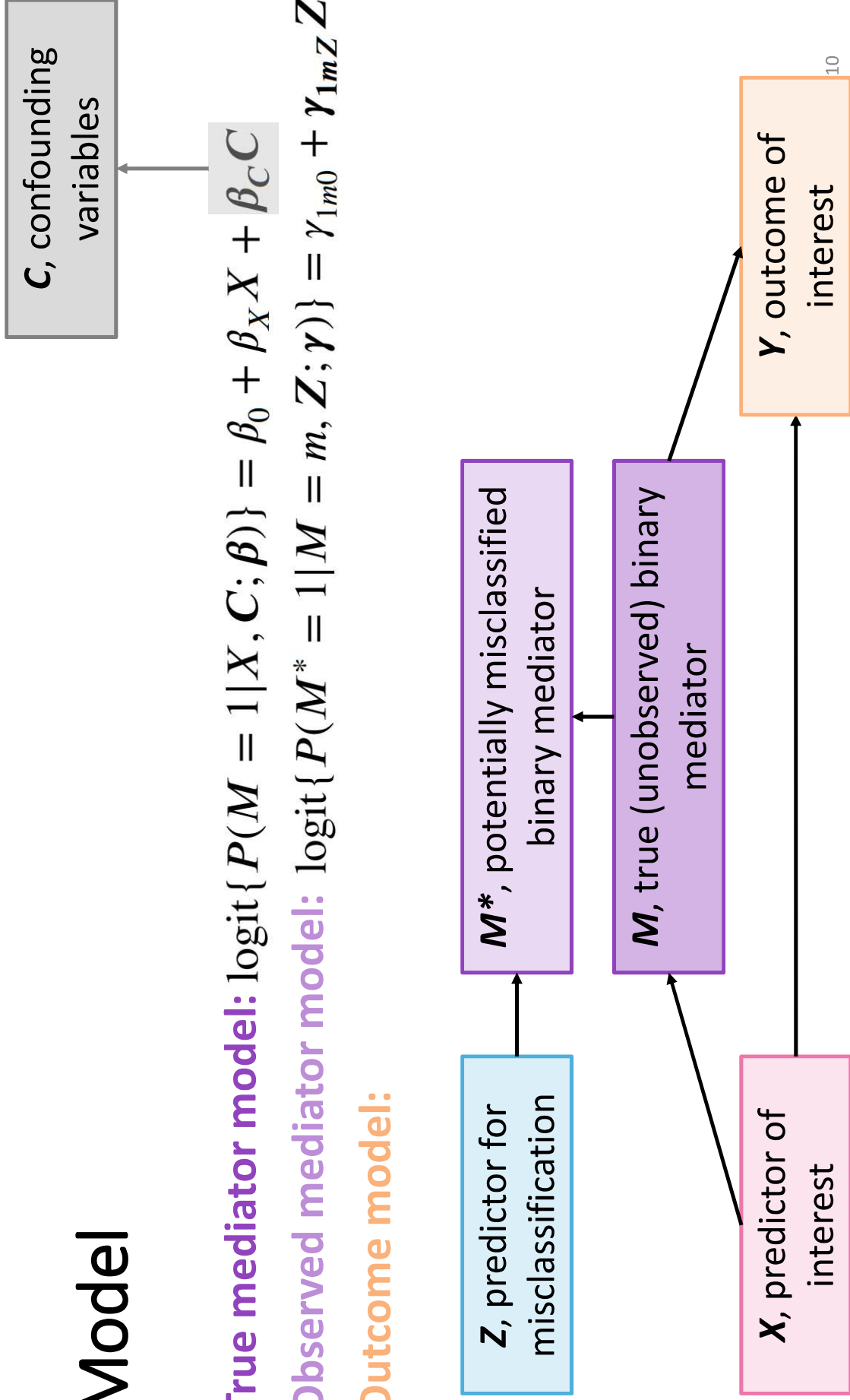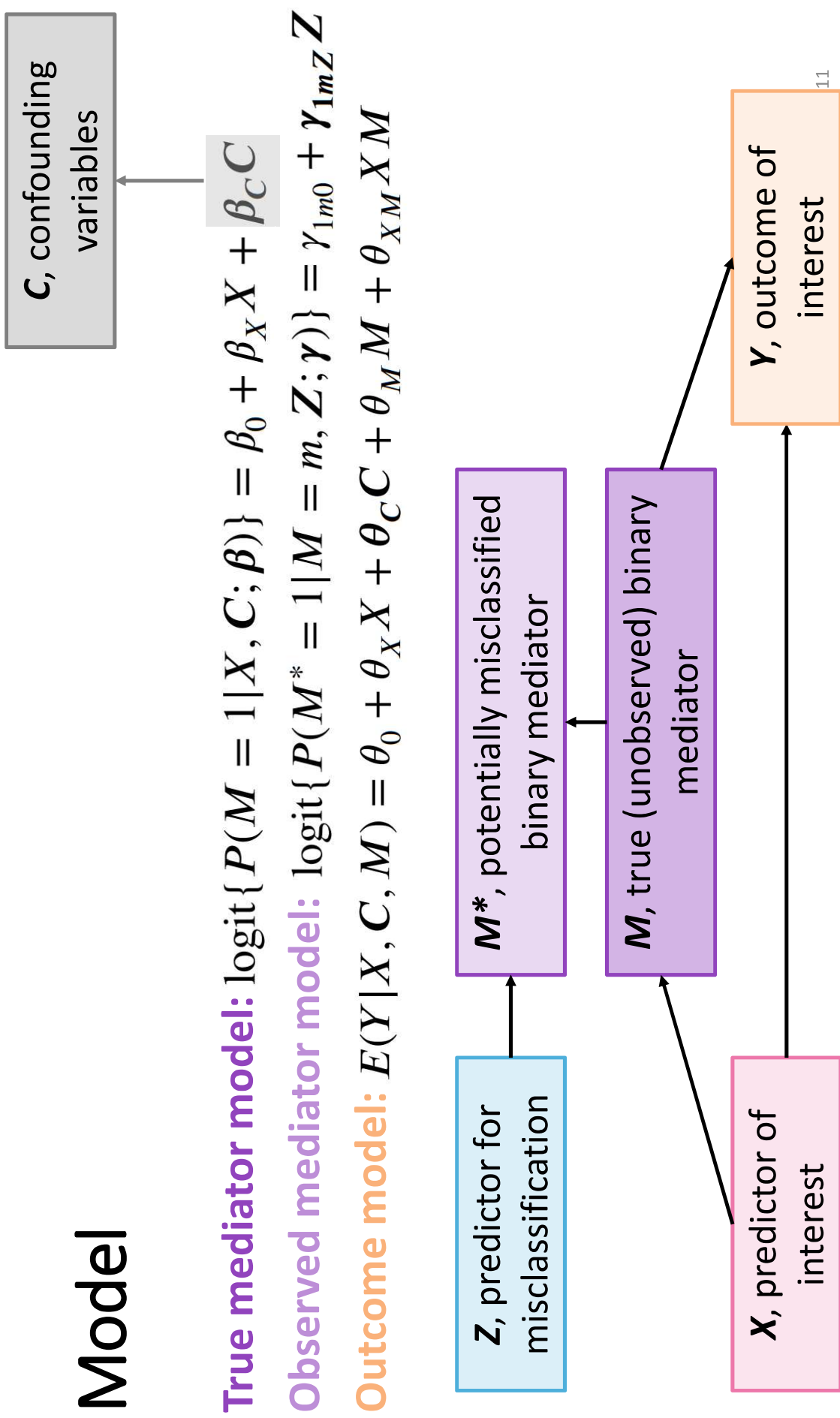
$X$, predictor of interest

# Model

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:**

**C**, confounding variables

**Z**, predictor for misclassification

**M***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**X**, predictor of interest

**Y**, outcome of interest

# Model

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$
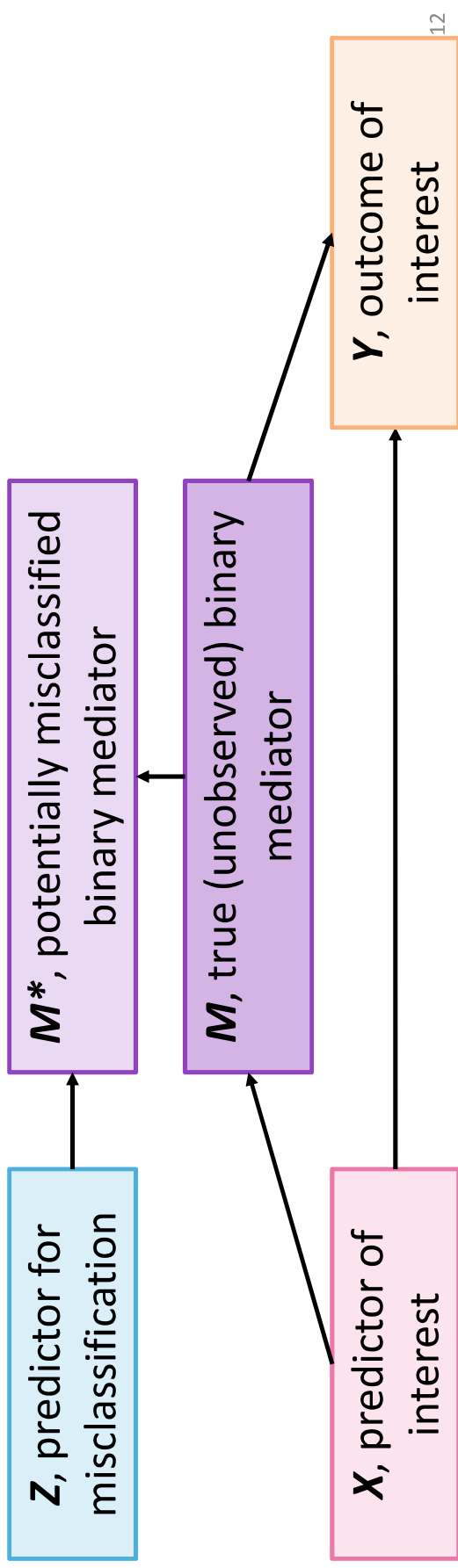
**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

**C**, confounding variables — $\beta_C C$

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**Y**, outcome of interest

**Z**, predictor for misclassification

**X**, predictor of interest

# Model

**Primary interest:** Estimating β and θ

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$
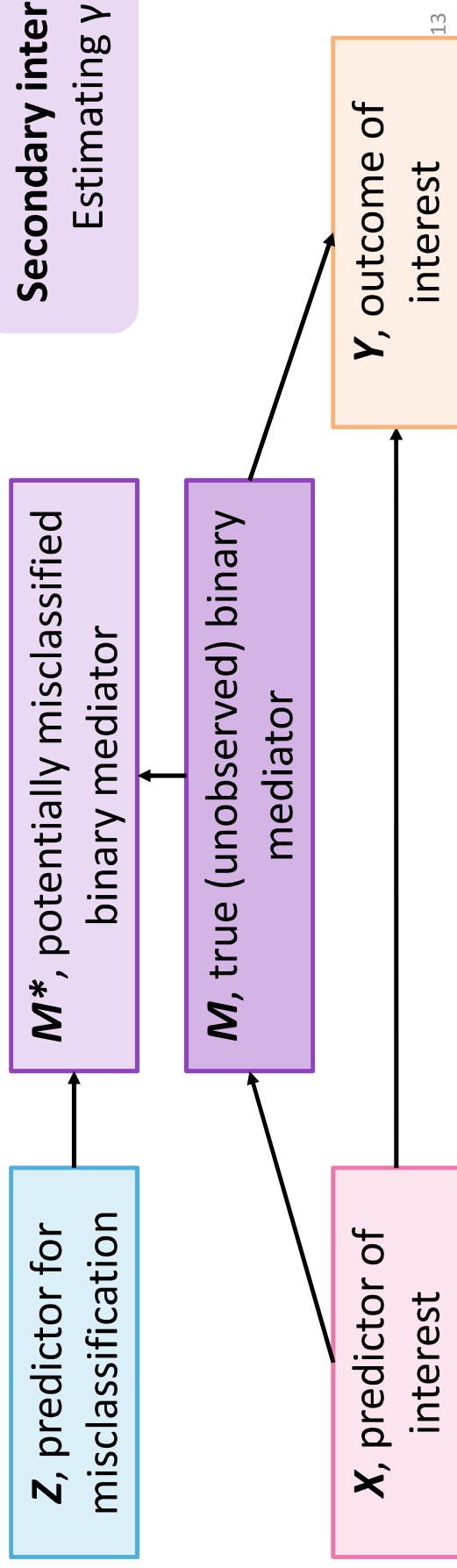
$M^*$, potentially misclassified binary mediator

$M$, true (unobserved) binary mediator

$Y$, outcome of interest

$Z$, predictor for misclassification

$X$, predictor of interest

# Model

**Primary interest:** Estimating β and θ

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

**Secondary interest:** Estimating γ



$M^*$, potentially misclassified binary mediator

$M$, true (unobserved) binary mediator

$Y$, outcome of interest

$Z$, predictor for misclassification

$X$, predictor of interest

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, \boldsymbol{C}; \boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_{\boldsymbol{C}} \boldsymbol{C}$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, \boldsymbol{Z}; \boldsymbol{\gamma})\} = \gamma_{1m0} + \boldsymbol{\gamma_{1mZ}} \boldsymbol{Z}$

**Outcome model:** $E(Y|X, \boldsymbol{C}, M) = \theta_0 + \theta_X X + \theta_{\boldsymbol{C}} \boldsymbol{C} + \theta_M M + \theta_{XM} XM$

| #1: OLS Correction | #2: Predictive value weighting | #3: An EM algorithm |
|---|---|---|

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, \boldsymbol{Z}; \boldsymbol{\gamma})\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, \boldsymbol{C}, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

| #1: OLS Correction[1] | #2: Predictive value weighting[2] | #3: An EM algorithm |

**Key point:** We can use **COMBO** to estimate subject-level sensitivity and specificity, and then plug these values into existing misclassification correction procedures.

- Existing procedures relied on *known* sensitivity and specificity.

1. Extended from Nguimkeu, Rosenman, and Tennekoon (2021), "Regression with a misclassified binary regressor: Correcting for hidden bias".
2. Extended from Lyles and Lin (2010), "Sensitivity analysis for misclassification in logistic regression via likelihood methods and PVW".

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, \boldsymbol{Z}; \boldsymbol{\gamma})\} = \gamma_{1m0} + \gamma_{1mZ} \boldsymbol{Z}$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#1: OLS Correction  #2: Predictive value weighting  #3: An EM algorithm

**Complete data log-likelihood:**

$$\ell_{complete}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\gamma}; X, C, \boldsymbol{Z}, Y)$$

$$= \sum_{i=1}^{N} \left[ \ell_{Y|X,M,C}(\boldsymbol{\theta}; X_i, M_i, C_i, Y_i) + \sum_{j=1}^{2} m_{ij} \log\{\pi_{ij}\} + \sum_{j=1}^{2} \sum_{\ell=1}^{2} m_{ij} m^*_{i\ell} \log\{\pi^*_{i\ell j}\} \right]$$

# Estimation

**True mediator model:** $\text{logit}\{P(M=1|X,C;\beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^*=1|M=m,Z;\gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X,C,M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

**Complete data log-likelihood:**

$$\ell_{complete}(\beta, \gamma, \gamma; X, C, Z, Y)$$

$$= \sum_{i=1}^{N} \left[ \ell_{Y|X,M,C}(\theta; X_i, M_i, C_i, Y_i) + \sum_{j=1}^{2} m_{ij}\log\{\pi_{ij}\} + \sum_{j=1}^{2}\sum_{\ell=1}^{2} m_{ij} m_{i\ell}^* \log\{\pi_{i\ell j}^*\} \right]$$

**Outcome**

# Estimation

**True mediator model:** $\operatorname{logit}\{P(M=1|X,C;\boldsymbol{\beta})\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\operatorname{logit}\{P(M^*=1|M=m,\boldsymbol{Z};\boldsymbol{\gamma})\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X,C,M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

**Complete data log-likelihood:**

$$\ell_{complete}(\boldsymbol{\beta},\boldsymbol{\gamma},\boldsymbol{\gamma};X,C,Z,Y)$$

$$= \sum_{i=1}^{N}\left[ \ell_{Y|X,M,C}(\boldsymbol{\theta};X_i,M_i,C_i,Y_i) + \sum_{j=1}^{2} m_{ij}\log\{\pi_{ij}\} + \sum_{j=1}^{2}\sum_{\ell=1}^{2} m_{ij}m_{i\ell}^{*}\log\{\pi_{i\ell j}^{*}\} \right]$$

**Outcome**

$I(M_i = j)$

$P(M_i = j)$

**True mediator**

18

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

**Complete data log-likelihood:**

$$\ell_{complete}(\beta, \Upsilon, \gamma; X, C, Z, Y)$$

$$= \sum_{i=1}^{N} \left[ \ell_{Y|X,M,C}(\theta; X_i, M_i, C_i, Y_i) + \sum_{j=1}^{2} m_{ij} \log\{\pi_{ij}\} + \sum_{j=1}^{2} \sum_{\ell=1}^{2} m_{ij} m^*_{i\ell} \log\{\pi^*_{i\ell j}\} \right]$$

**Outcome**

**True mediator**

**Observed mediator**

$I(M_i = j)$    $P(M_i = j)$

$I(M_i^* = \ell)$    $P(M_i^* = \ell \mid M_i = j)$

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#1: OLS Correction

#2: Predictive value weighting

#3: An EM algorithm

**Expectation Step**

**Maximization Step**

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1 | X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1 | M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y | X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#1: OLS Correction

#2: Predictive value weighting

#3: An EM algorithm

**Expectation Step**

**Maximization Step**

$$w_{ij} = P(M_i = j | M_i^*, X_i, C_i, Z_i, Y_i)$$

$$= \sum_{\ell=1}^{2} \frac{m_{i\ell}^* \pi_{i\ell j}^* \pi_{ij} E[Y_i | X_i, M_i = j, C_i, \theta^{(t)}]}{\sum_{k=1}^{2} \pi_{i\ell k}^* \pi_{ik} E[Y_i | X_i, M_i = k, C_i, \theta^{(t)}]}$$

# Estimation

**True mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#1: OLS Correction

#2: Predictive value weighting

#3: An EM algorithm

**Expectation Step**

$$w_{ij} = P(M_i = j|M_i^*, X_i, C_i, Z_i, Y_i)$$

$$= \sum_{\ell=1}^{2} \frac{m_{i\ell j}^* \pi_{i\ell j}^* \pi_{ij} E[Y_i|X_i, M_i = j, C_i, \theta^{(t)}]}{\sum_{k=1}^{2} \pi_{i\ell k}^* \pi_{ik} E[Y_i|X_i, M_i = k, C_i, \theta^{(t)}]}$$

**Maximization Step**

$$Q = \sum_{i=1}^{N} \sum_{j=1}^{2} \left[ \sum_{j=1}^{2} \ell_{Y|X,M,C}(\theta; X_i, M_i = w_{ij}, C_i, Y_i) \right.$$

$$\left. + \sum_{j=1}^{2} w_{ij} \log\{\pi_{ij}\} + \sum_{j=1}^{2} \sum_{\ell=1}^{2} w_{ij} m_{i\ell}^* \log\{\pi_{i\ell j}^*\} \right]$$

# Estimation

**Mediator model:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

$\text{logit}\{P(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#3: An EM algorithm

**Maximization Step**

#2: Predictive value weighting

#1: OLS Correction

$$Q_\beta = \sum_{i=1}^{N} \Big[ \sum_{j=1}^{2} w_{ij}\log\{\pi_{ij}\} \Big]$$

$$Q_{\gamma_1} = \sum_{i=1}^{N} \Big[ \sum_{\ell=1}^{2} w_{i1} m^*_{i\ell}\log\{\pi^*_{i\ell 1}\} \Big]$$

$$Q_{\gamma_2} = \sum_{i=1}^{N} \Big[ \sum_{\ell=1}^{2} w_{12} m^*_{i\ell}\log\{\pi^*_{i\ell 2}\} \Big]$$

$$_u Q_\theta = \sum_{i=1}^{N} \Big[ \sum_{j=1}^{2} \ell_{Y|X,M,C}(\theta; X_i, M_i = w_{ij}, C_i, Y_i) \Big]$$

$$= \sum_{\ell=1}^{2} \frac{m^*_{i\ell} \pi^*_{i\ell j} \pi_{ij} E[Y_i|X_i, M_i = j, C_i, \theta^{(t)}]}{\sum_{k=1}^{2} \pi^*_{i\ell k} \pi_{ik} E[Y_i|X_i, M_i = k, C_i, \theta^{(t)}]}$$

$$Q = \sum_{i=1}^{N} \Big[ \sum_{j=1}^{2} \ell_{Y|X,M,C}(\theta; X_i, M_i = w_{ij}, C_i, Y_i) \Big]$$
$$+ \sum_{j=1}^{2} w_{ij}\log\{\pi_{ij}\} + \sum_{j=1}^{2}\sum_{\ell=1}^{2} w_{ij} m^*_{i\ell}\log\{\pi^*_{i\ell j}\}$$

# Estimation

**True mediator model:** $\text{logit}\{P(M=1|X,C;\beta)\} = \beta_0 + \beta_X X + \beta_C C$

**Observed mediator model:** $\text{logit}\{P(M^*=1|M=m,Z;\gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$

**Outcome model:** $E(Y|X,C,M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} XM$

#1: OLS Correction  #2: Predictive value weighting  #3: An EM algorithm

Apply the label switching correction from Webb and Wells (2023)

**Expectation Step**

$$w_{ij} = P(M_i = j | M_i^*, X_i, C_i, Z_i, Y_i)$$

$$= \sum_{\ell=1}^{2} \frac{m_{i\ell}^* \pi_{i\ell j}^* \pi_{ij} E[Y_i | X_i, M_i = j, C_i, \theta^{(t)}]}{\sum_{k=1}^{2} \pi_{i\ell k}^* \pi_{ik} E[Y_i | X_i, M_i = k, C_i, \theta^{(t)}]}$$

**Maximization Step**

$$Q = \sum_{i=1}^{N} \sum_{j=1}^{2} \Big[ \sum_{j=1}^{2} \ell_{Y|X,M,C}(\theta; X_i, M_i = w_{ij}, C_i, Y_i)$$

$$+ \sum_{j=1}^{2} w_{ij} \log\{\pi_{ij}\} + \sum_{j=1}^{2} \sum_{\ell=1}^{2} w_{ij} m_{i\ell}^* \log\{\pi_{i\ell j}^*\} \Big]$$

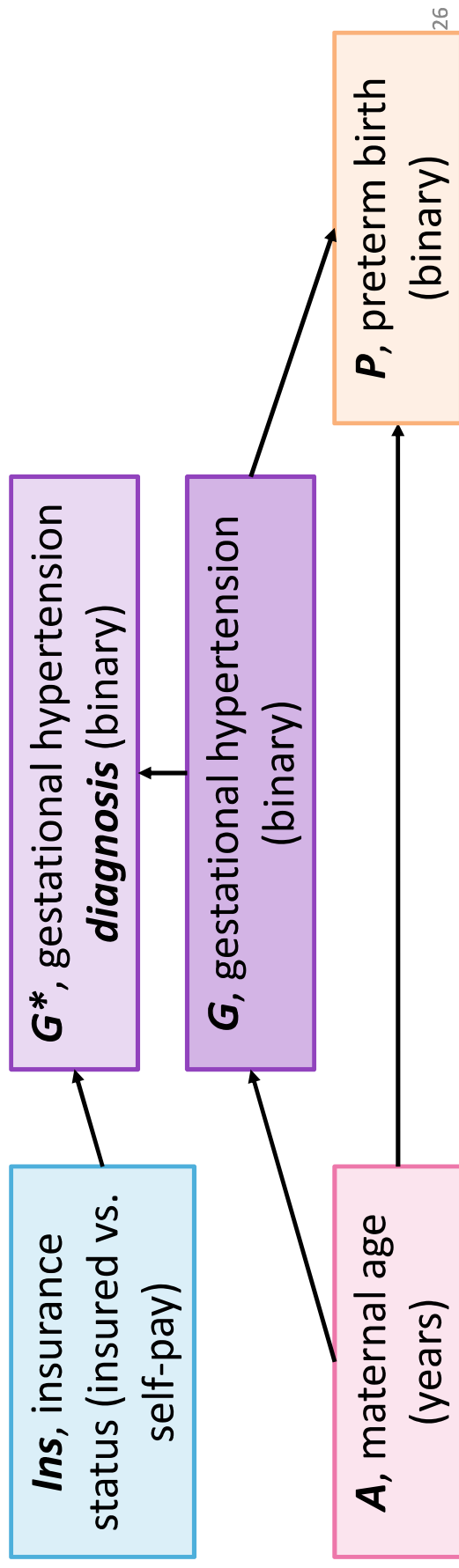24

# Problem setting

- **Example:**
  - Does **gestational hypertension** mediate the association between **maternal age** and **preterm birth**, after accounting for potential **misdiagnosis of gestational hypertension** based on **patient insurance status**?

**Ins**, insurance status (insured vs. self-pay)

**G\***, gestational hypertension **diagnosis** (binary)

**G**, gestational hypertension (binary)

**A**, maternal age (years)

**P**, preterm birth (binary)

# Applied Example

**Data:** National Vital Statistics System of the National Center for Health Statistics
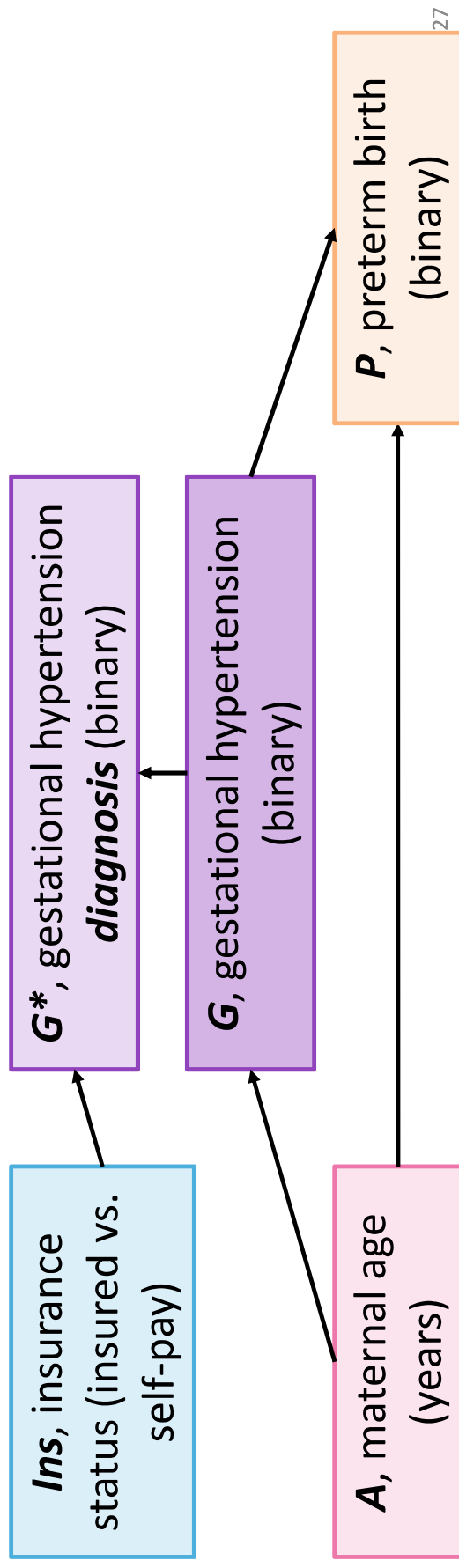
- Random sample of 20,000 observations.



**Ins**, insurance status (insured vs. self-pay)

**G\***, gestational hypertension *diagnosis* (binary)

**G**, gestational hypertension (binary)

**A**, maternal age (years)

**P**, preterm birth (binary)

# Applied Example

**True mediator model:** $G \sim$ Age + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** $G^* \mid G \sim$ Race + Ins

**Outcome model:** $P \sim$ Age + Race + Education + Parity + Smoking + BMI + $G$ + $G$ * Age



*Ins*, insurance status (insured vs. self-pay)

*A*, maternal age (years)

*G\**, gestational hypertension *diagnosis* (binary)

*G*, gestational hypertension (binary)

*P*, preterm birth (binary)

# Applied Example

**True mediator model:** **G** ~ Age + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** **G\*** | **G** ~ Race + Ins

**Outcome model:** **P** ~ Age + Race + Education + Parity + Smoking + BMI + **G** + **G \* Age**

| | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
| | Est. | SE | Est. | SE |
| $\beta_{age}$ | | | | |
| $\gamma_{ins, G = 1}$ | | | | |
| $\gamma_{ins, G = 2}$ | | | | |
| $\theta_{age}$ | | | | |
| $\theta_{G}$ | | | | |
| $\theta_{G * age}$ | | | | |

# Applied Example

**True mediator model:** **G** ~ <mark>Age</mark> + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** **G\*** | **G** ~ Race + <mark>Ins</mark>

**Outcome model:** <mark>**P** ~ Age + Race + Education + Parity + Smoking + BMI + **G** + **G** \* Age</mark>

| | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
| | Est. | SE | Est. | SE |
| $\beta_{age}$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{ins,\,G=1}$ | | | | |
| $\gamma_{ins,\,G=2}$ | | | | |
| $\theta_{age}$ | | | | |
| $\theta_{G}$ | | | | |
| $\theta_{G\,*\,age}$ | | | | |

Association between **age** and **G** unchanged, accounting for misdiagnosis

# Applied Example

**True mediator model:** $G$ ~ Age + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** $G^* \mid G$ ~ Race + Ins

**Outcome model:** $P$ ~ Age + Race + Education + Parity + Smoking + BMI + $G$ + $G$ * Age

| | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
| | Est. | SE | Est. | SE |
| $\beta_{age}$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{ins,\,G=1}$ | | | | |
| $\gamma_{ins,\,G=2}$ | | | | |
| $\theta_{age}$ | 0.02 | 0.05 | 0.10 | 0.03 |
| $\theta_{G}$ | 1.19 | 0.17 | 0.88 | 0.06 |
| $\theta_{G\,*\,age}$ | 0.19 | 0.09 | 0.06 | 0.06 |

Association between **age** and **G** unchanged, accounting for misdiagnosis

Association between **G** and **P** strengthens

# Applied Example

**True mediator model:** $G \sim$ **Age** + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** **G\*** | **G** $\sim$ Race + Ins

**Outcome model:** **P** $\sim$ Age + Race + Education + Parity + Smoking + BMI + **G** + **G \* Age**

Use γ estimates to compute sensitivity and specificity.

|  | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
|  | Est. | SE | Est. | SE |
| $\beta_{age}$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{ins, G=1}$ | -1.01 | 0.40 | - | - |
| $\gamma_{ins, G=2}$ | 2.09 | 8.81 | - | - |
| $\theta_{age}$ | 0.02 | 0.05 | 0.10 | 0.03 |
| $\theta_G$ | 1.19 | 0.17 | 0.88 | 0.06 |
| $\theta_{G * age}$ | 0.19 | 0.09 | 0.06 | 0.06 |

Association between **age** and **G** unchanged, accounting for misdiagnosis

Association between **G** and **P** strengthens

# Applied Example

**True mediator model:** **G** ~ Age + Race + Education + Parity + Smoking + BMI

**Observed mediator model:** **G\*** | **G** ~ Race + Ins

**Outcome model:** **P** ~ Age + Race + Education + Parity + Smoking + BMI + **G** + **G** \* Age

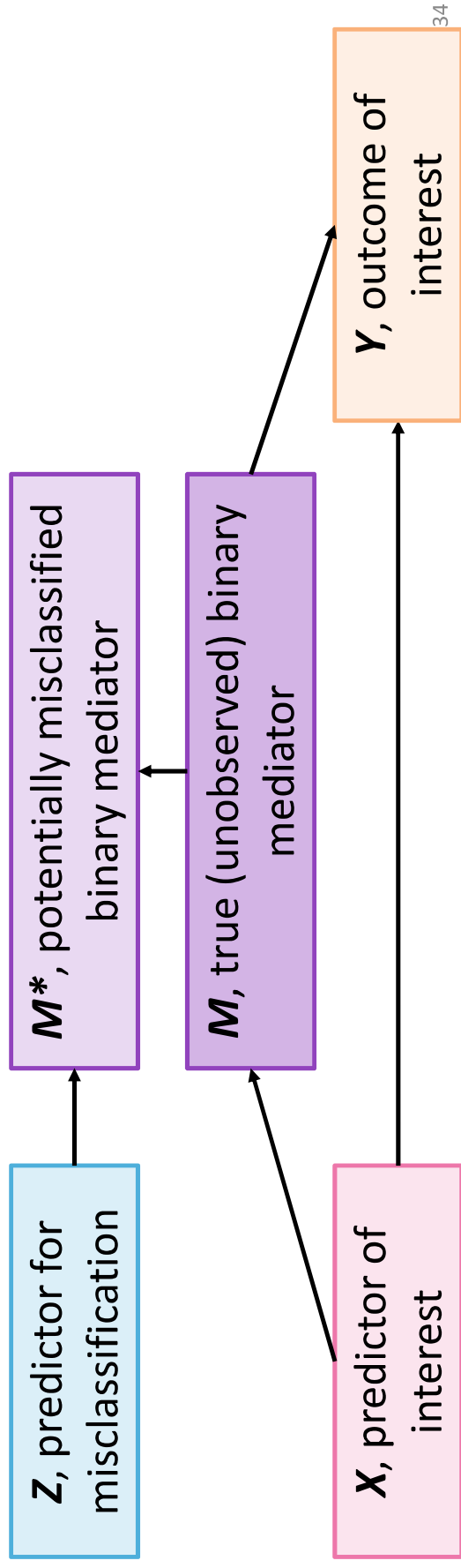| | Estimated Specificity<br>P( no G\* \| no G ) | Estimated Sensitivity<br>P( G\* \| G ) |
|---|---|---|
| Insured | 99.9% | 43.1% |
| Self-Pay | 99.4% | 21.7% |

# Code is available!

- Sample function and simulation code at:

**bit.ly/enar2024-code-webb**

# Conclusions and Next Steps

- We can use the proposed methods to **estimate associations** when a **binary mediator is potentially misclassified**.



**Z**, predictor for misclassification

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

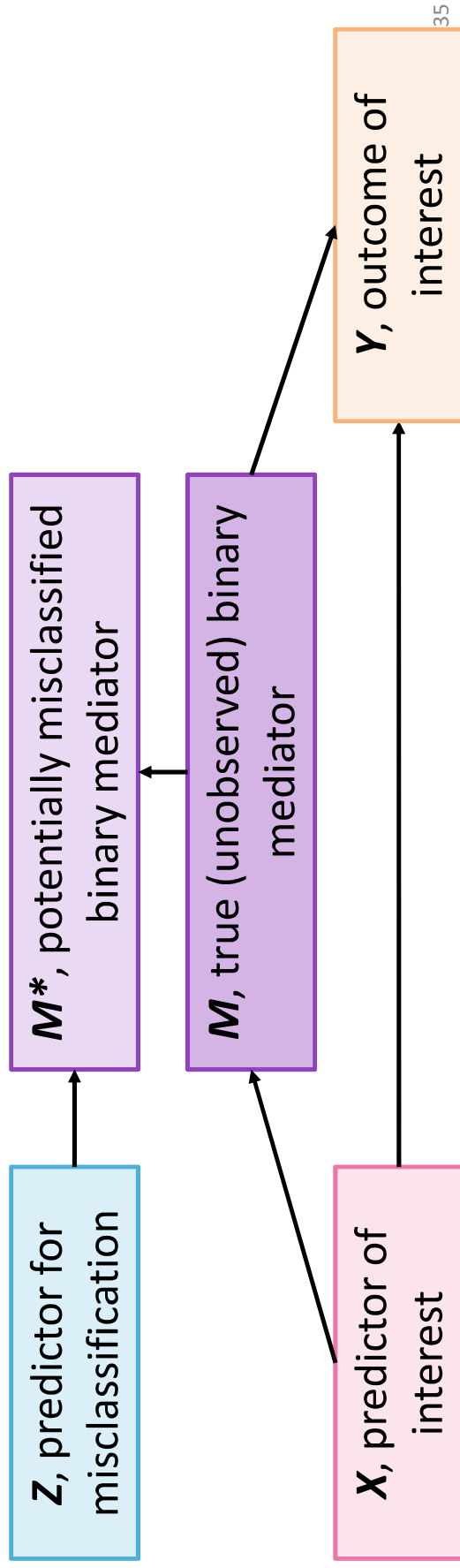**X**, predictor of interest

**Y**, outcome of interest

# Conclusions and Next Steps

- We can use the proposed methods to **estimate associations** when a **binary mediator is potentially misclassified.**

- **Next steps:** Incorporate other variables measured with error.



**Z**, predictor for misclassification

**M\***, potentially misclassified binary mediator

**M**, true (unobserved) binary mediator

**X**, predictor of interest

**Y**, outcome of interest

# Thank you!

## Kimberly A. H. Webb

kah343@cornell.edu

kimhwebb.com —→ My "webb-site" ☺

Cornell Bowers C·IS
**Statistics and Data Science**