# The misdiagnosed mediator:
## Estimating the effect of maternal age on preterm birth risk in the presence of misclassified gestational hypertension

**Kimberly A. H. Webb** and Martin T. Wells

Joint Statistical Meetings

August 5, 2025

# Mediation analysis

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

# Mediation analysis

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

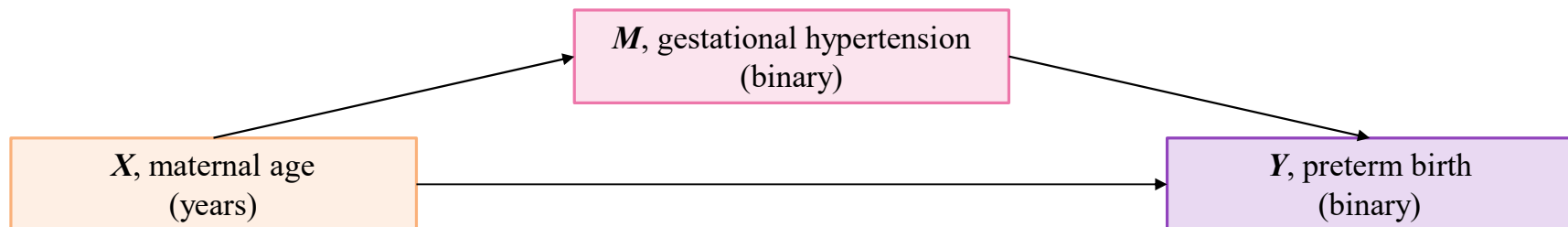| $X$, maternal age (years) | → | $Y$, preterm birth (binary) |

# Mediation analysis

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

# Mediation analysis

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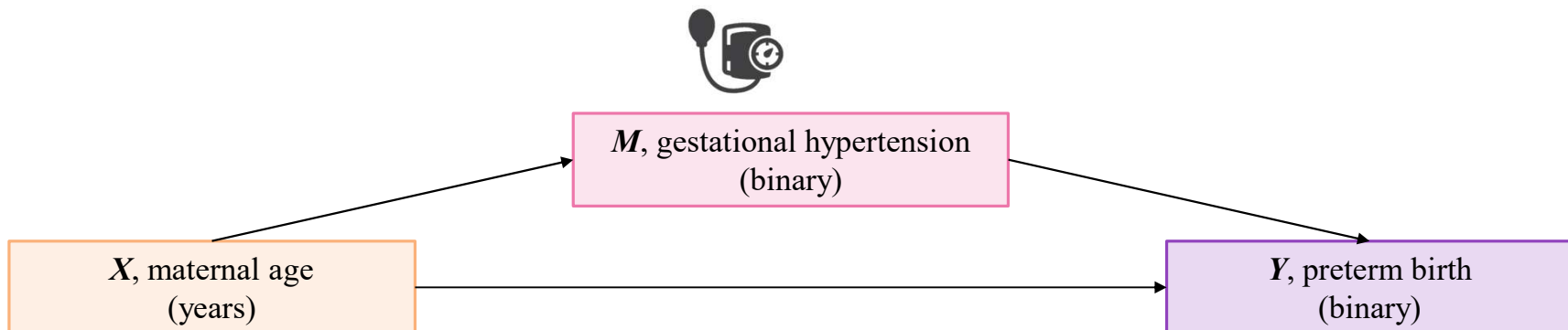
# Mediation analysis

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

Measurement of **M** is **not always accurate**, and we obtain **M***.

*M*, gestational hypertension (binary)

*X*, maternal age (years)
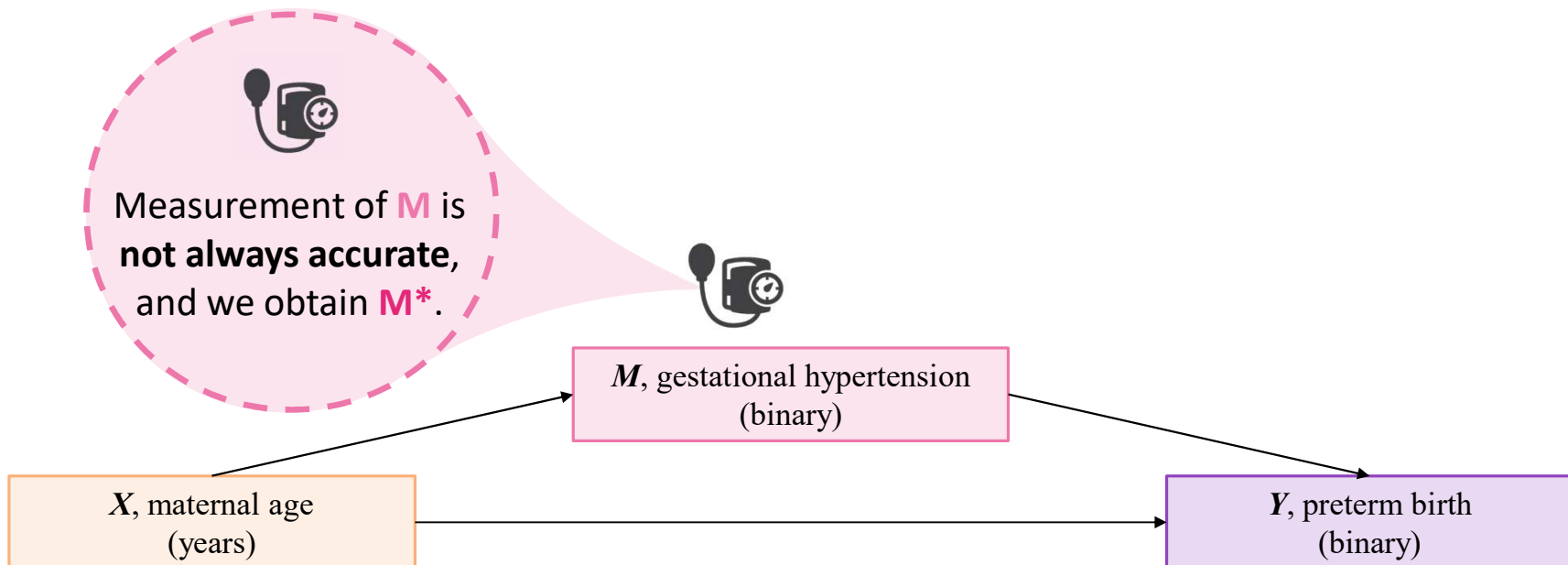
*Y*, preterm birth (binary)

# Mediation analysis

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



Measurement of **M** is **not always accurate**, and we obtain **M***.

*M**, gestational hypertension *diagnosis* (binary)

*M*, gestational hypertension (binary)

*X*, maternal age (years)

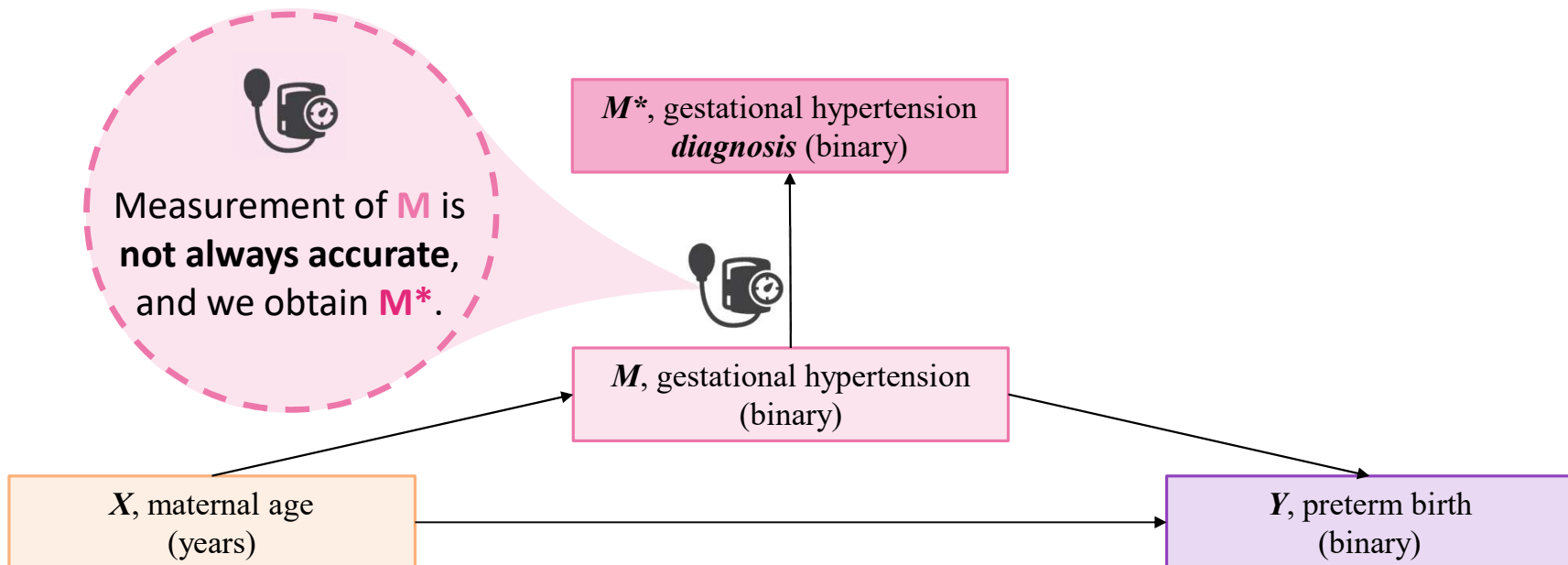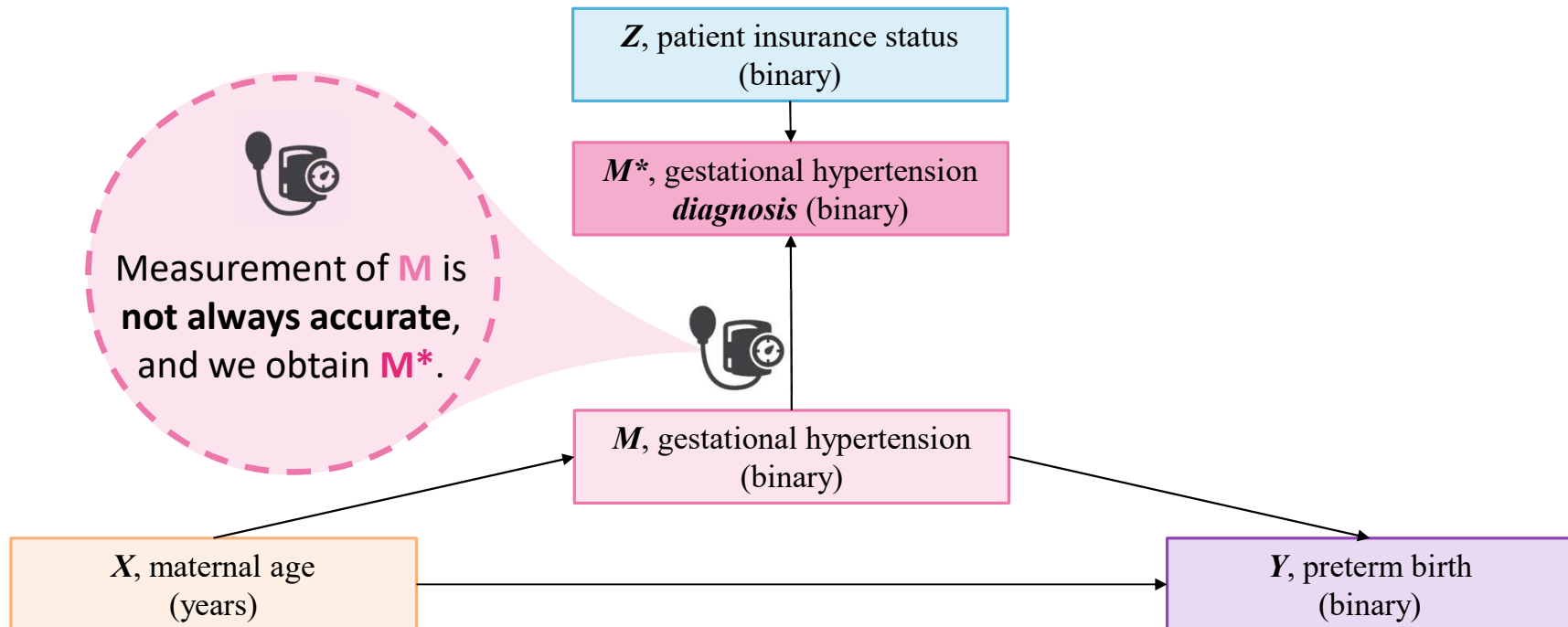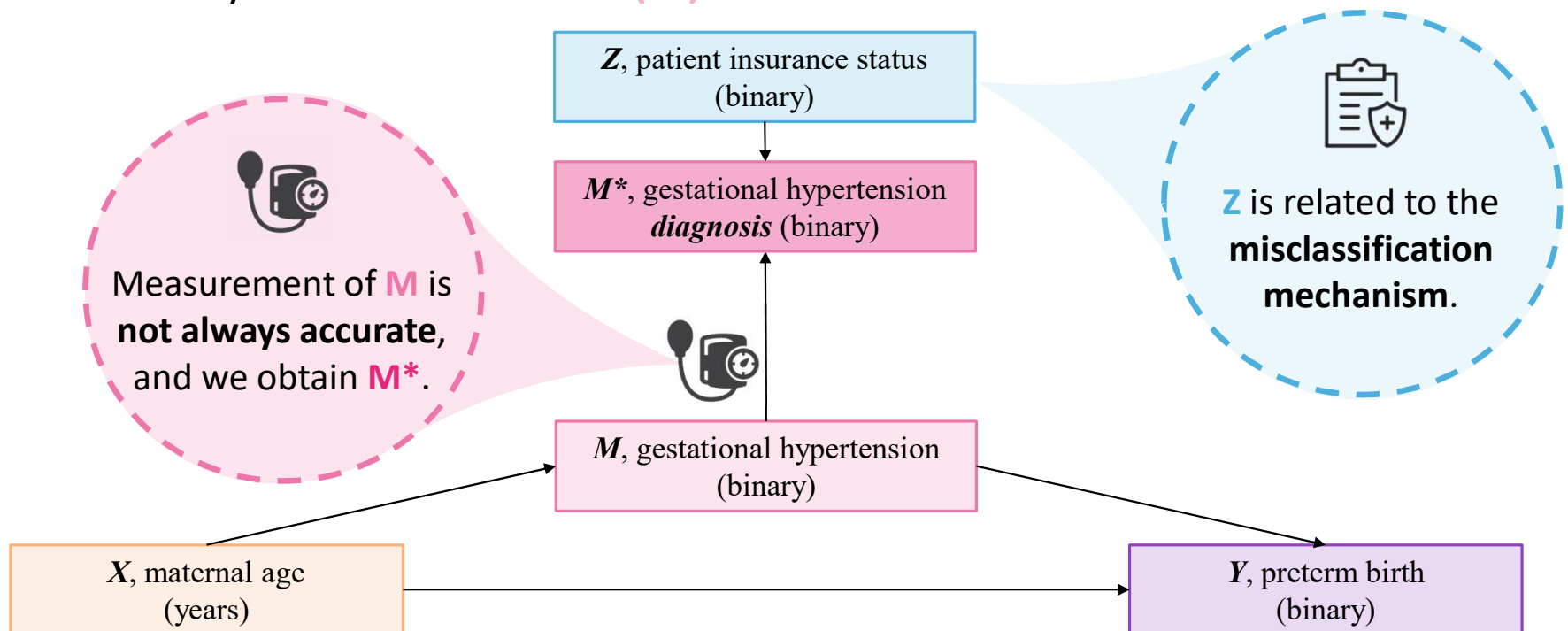*Y*, preterm birth (binary)

# Mediation analysis

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

# Mediation analysis

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



Measurement of **M** is **not always accurate**, and we obtain **M\***.

**Z**, patient insurance status (binary)

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

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

**M**, gestational hypertension (binary)

**X**, maternal age (years)

**Y**, preterm birth (binary)

# Our analysis is complicated by misdiagnosis in *M*

# Our analysis is complicated by misdiagnosis in *M*



$Z$, patient insurance status (binary)

$M^*$, gestational hypertension *diagnosis* (binary)

Misclassification is **covariate-dependent**.

$M$, gestational hypertension (binary)

$X$, maternal age (years)

$Y$, preterm birth (binary)

# Our analysis is complicated by misdiagnosis in *M*

No **gold standard** measure for **M**.
- M is **latent**.

*Z*, patient insurance status (binary)

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

Misclassification is **covariate-dependent**.

*M*, gestational hypertension (**latent**, binary)

*X*, maternal age (years)

*Y*, preterm birth (binary)

# Our analysis is complicated by misdiagnosis in *M*

Ignoring misclassification in M* → Bias in parameter and effect estimates

No **gold standard** measure for **M**.
- M is **latent**.

Misclassification is **covariate-dependent**.

**Z**, patient insurance status (binary)

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

**M**, gestational hypertension (**latent**, binary)

**X**, maternal age (years)

**Y**, preterm birth (binary)

# Analysis Plan



Define the **misclassification model**

**1**

Develop **estimation methods** for parameters of interest

**2**

Use parameter estimates to compute **(in)direct effects and misclassification rates**

**3**

# Our model accounts for misdiagnosis in *M*

# Our model accounts for misdiagnosis in *M*

**Z**, patient insurance status (binary)

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

True mediator mechanism

**M**, gestational hypertension (**latent**, binary)

**X**, maternal age (years)

**Y**, preterm birth (binary)

# Our model accounts for misdiagnosis in *M*



Z, patient insurance status (binary)

M*, gestational hypertension *diagnosis* (binary)

True mediator mechanism

*M*, gestational hypertension (**latent**, binary)

*X*, maternal age (years)

Y, preterm birth (binary)

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

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



$C$, covariates

$Z$, patient insurance status (binary)

$M^*$, gestational hypertension *diagnosis* (binary)

True mediator mechanism

$M$, gestational hypertension (**latent**, binary)

$X$, maternal age (years)

$Y$, preterm birth (binary)

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

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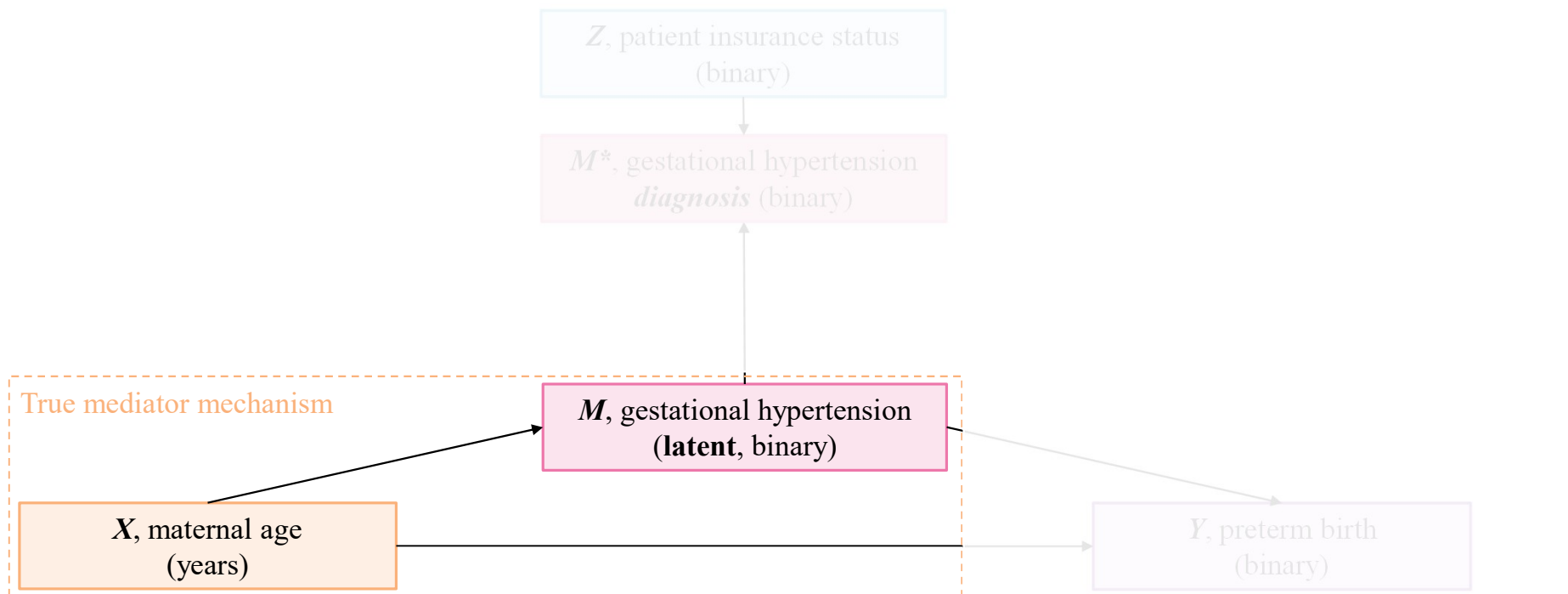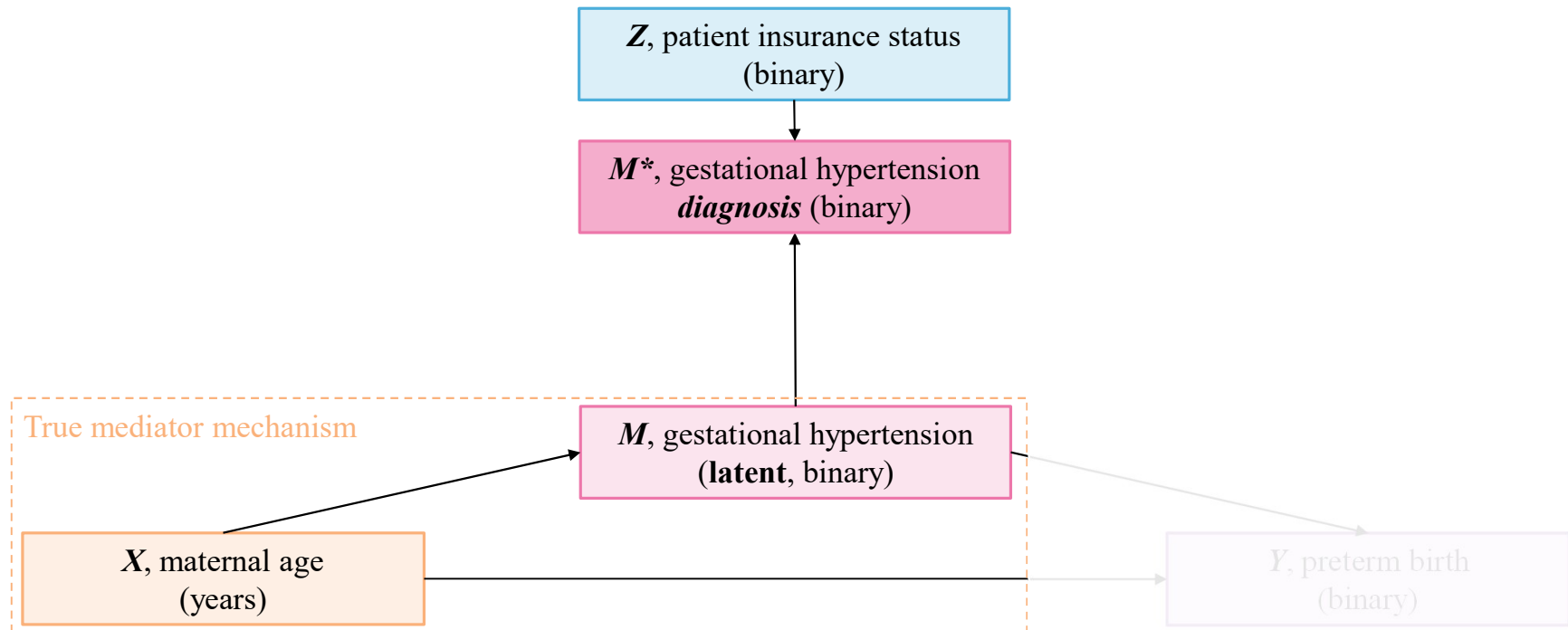
**True mediator mechanism:** $\text{logit}\{P(M = 1|X, C; \beta)\} = \beta_0 + \beta_X X + \beta_C C$ ①

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

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

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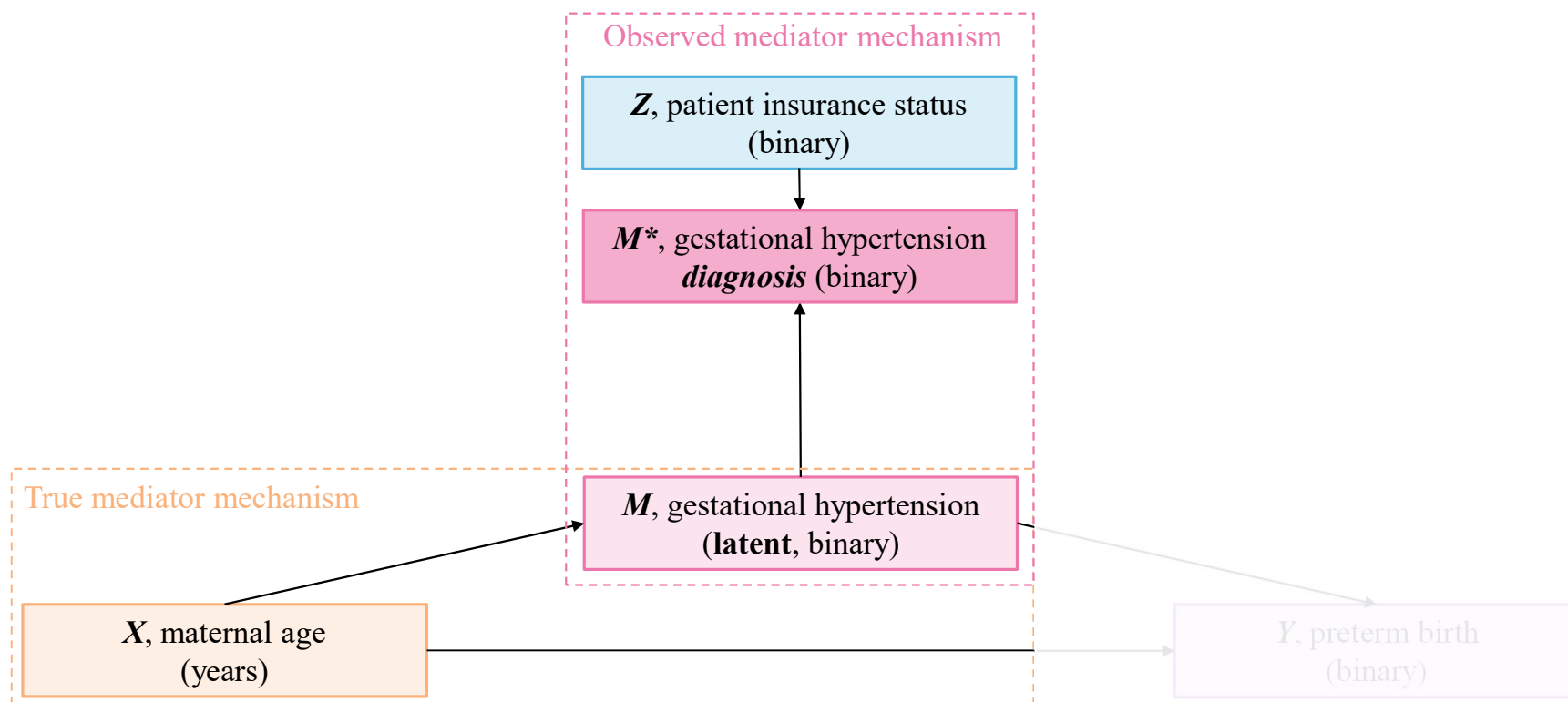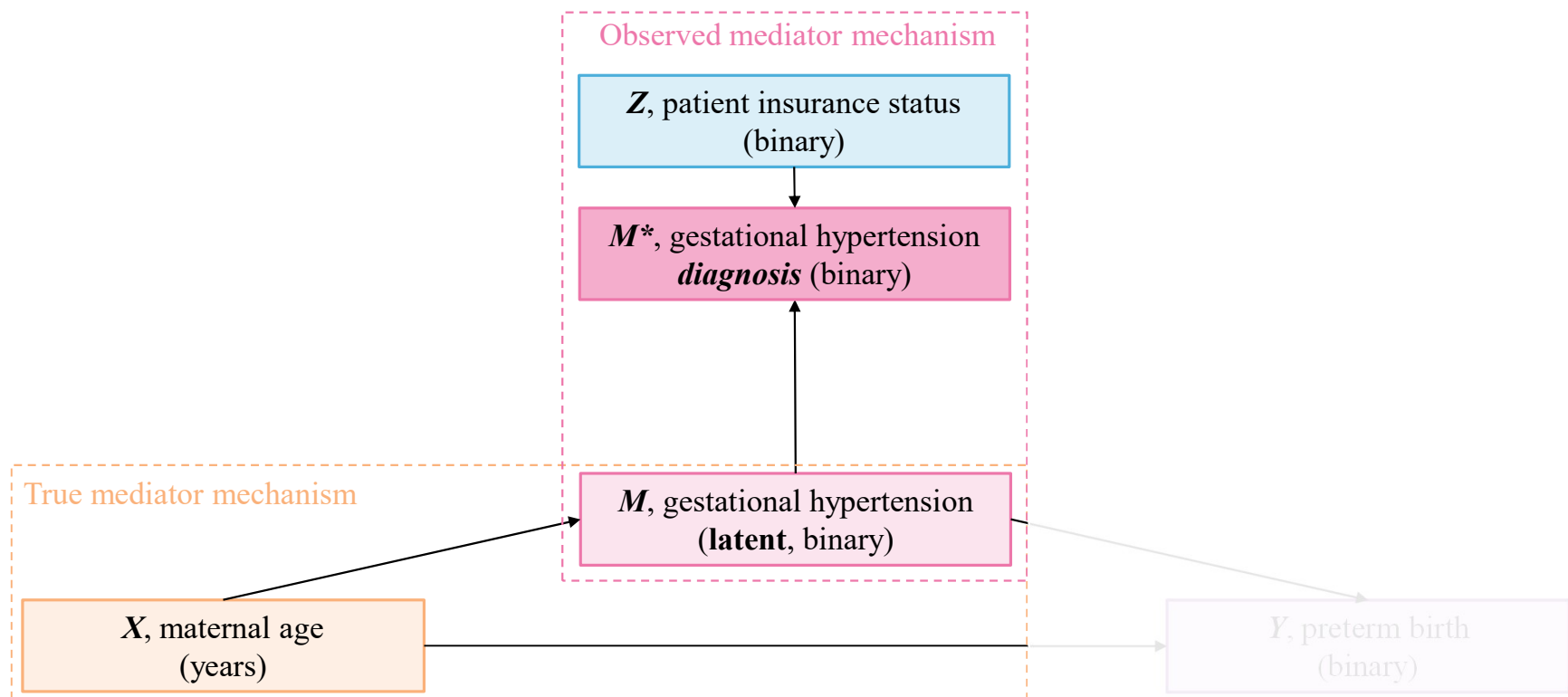
①

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

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

①



Observed mediator mechanism

$Z$, patient insurance status (binary)

$M^*$, gestational hypertension *diagnosis* (binary)

True mediator mechanism

$M$, gestational hypertension (**latent**, binary)

Outcome mechanism

$X$, maternal age (years)

$Y$, preterm birth (binary)

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

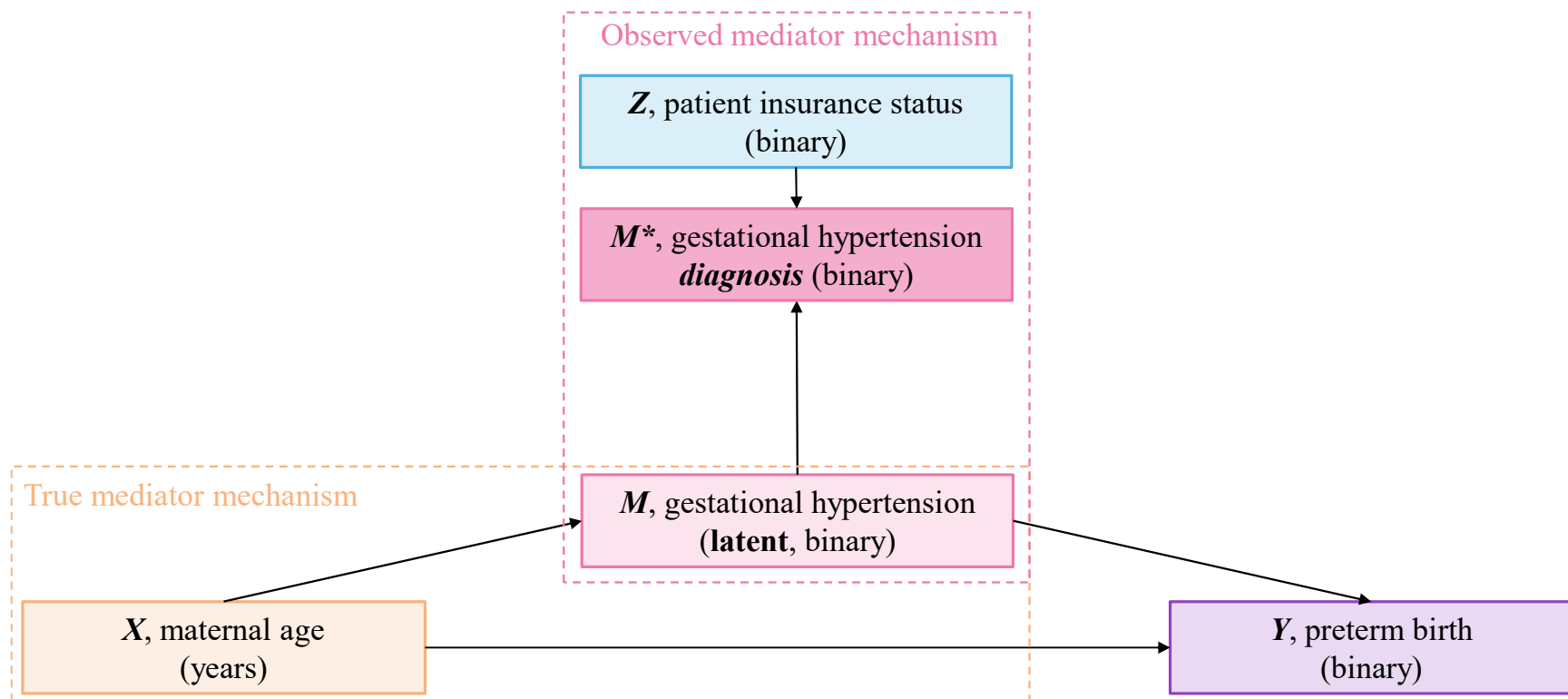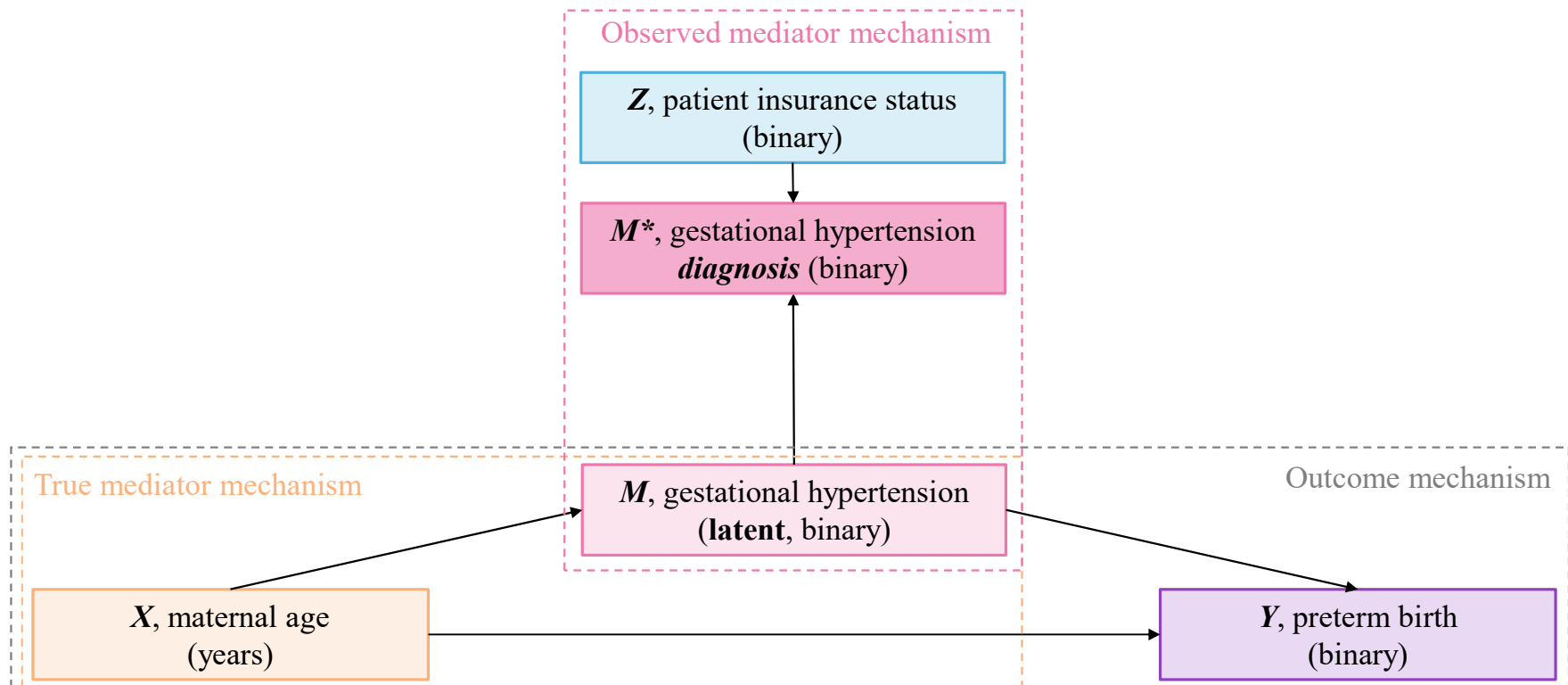**Observed mediator mechanism:** $\text{logit}\{P(M^* = 1|M = m, Z; \gamma)\} = \gamma_{1m0} + \gamma_{1mZ} Z$
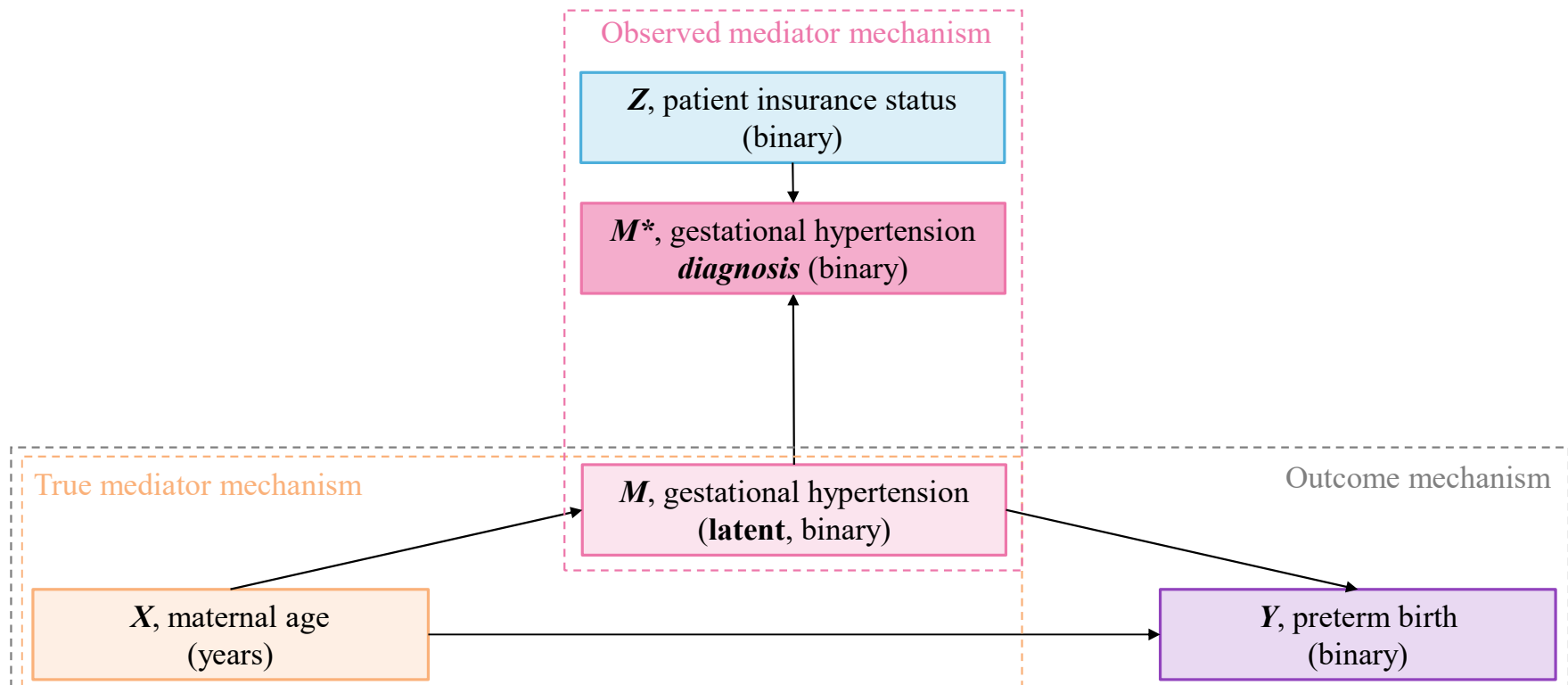
**Outcome mechanism:** $E(Y|X, C, M) = \theta_0 + \theta_X X + \theta_C C + \theta_M M + \theta_{XM} X M$

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

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

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



**Mediation parameters:**
β and θ quantify the exposure-mediator-outcome relationship.

Observed mediator mechanism

$Z$, patient insurance status (binary)

$M^*$, gestational hypertension *diagnosis* (binary)

True mediator mechanism

$M$, gestational hypertension (**latent**, binary)

Outcome mechanism

$X$, maternal age (years)

$Y$, preterm birth (binary)

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

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

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

1

**Mediation parameters:** β and θ quantify the exposure-mediator-outcome relationship.

Observed mediator mechanism

$\boldsymbol{Z}$, patient insurance status (binary)

$\boldsymbol{M^*}$, gestational hypertension *diagnosis* (binary)

**Misclassification parameters:** γ quantifies the effect of Z on misclassification rates

True mediator mechanism

$\boldsymbol{M}$, gestational hypertension (**latent**, binary)

Outcome mechanism

$\boldsymbol{X}$, maternal age (years)

$\boldsymbol{Y}$, preterm birth (binary)

# We developed 3 estimation methods

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

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

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

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

# We developed 3 estimation methods

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

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

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

| #1: OLS Correction | #2: Predictive value weighting | #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".

# We developed 3 estimation methods

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

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

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

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

**Complete data log-likelihood:**

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

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

# We developed 3 estimation methods

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

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

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

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

**Complete data log-likelihood:**

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

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

Outcome

# We developed 3 estimation methods

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

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

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

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

**Complete data log-likelihood:**

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

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

$P(M_i = j)$

$\mathbb{I}(M_i = j)$

Outcome          True mediator

# We developed 3 estimation methods

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

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

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

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

**Complete data log-likelihood:**

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

$P(M_i = j)$

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

$\mathbb{I}(M_i = j)$

$\mathbb{I}(M_i^* = \ell)$

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

Outcome     **True mediator**     **Observed mediator**

# We developed 3 estimation methods

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

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

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

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

**Expectation Step** ⟵⟶ **Maximization Step**

# We developed 3 estimation methods

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

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

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

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

**Expectation Step** ⟷ **Maximization Step**

$$w_{ij} = P(M_i = j|M_i^*, X_i, \boldsymbol{C}_i, \boldsymbol{Z}_i, Y_i)$$

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

# We developed 3 estimation methods

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

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

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

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

**Expectation Step**

**Maximization Step**

$$w_{ij} = P(M_i = j|M_i^*, X_i, \boldsymbol{C}_i, \boldsymbol{Z}_i, Y_i)$$

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

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

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

# We developed 3 estimation methods

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

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

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

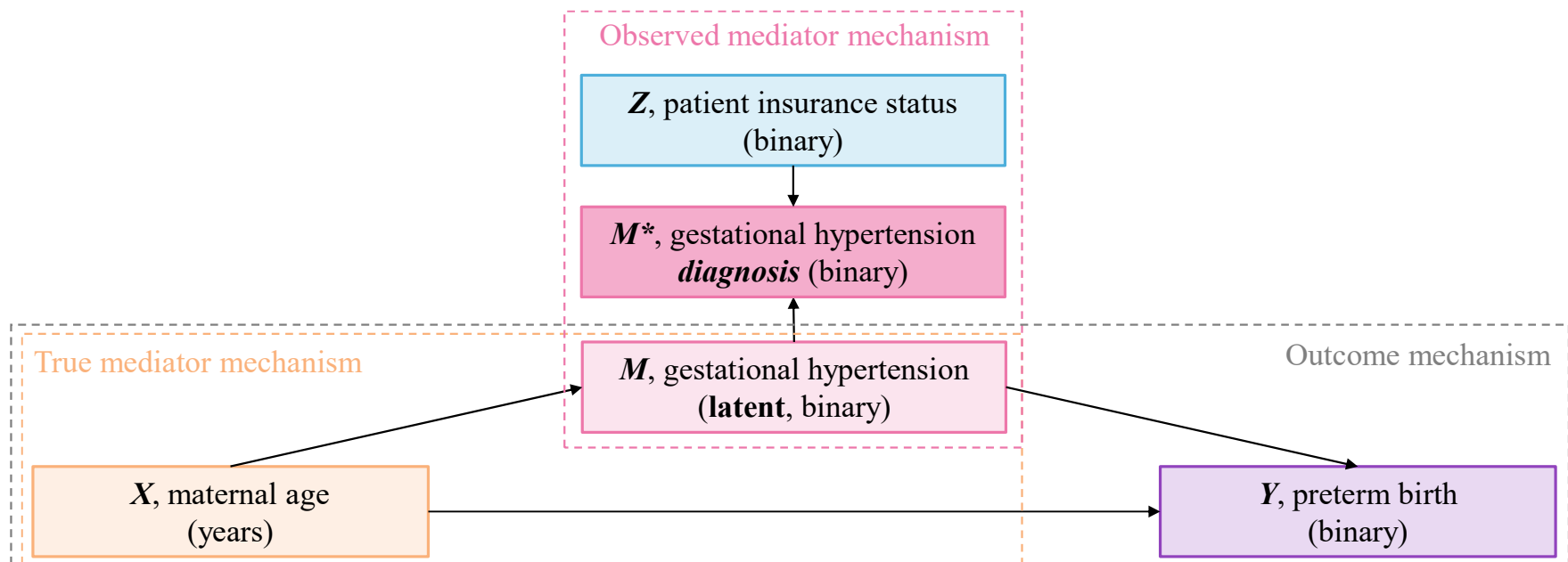| #1: OLS Correction | #2: Predictive value weighting | #3: An EM algorithm |
|---|---|---|

**Expectation Step**

Apply label switching correction
from Webb and Wells (2023)

**Maximization Step**

$$w_{ij} = P(M_i = j|M_i^*, X_i, \boldsymbol{C}_i, \boldsymbol{Z}_i, Y_i)$$

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

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

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

# We developed 3 estimation methods

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

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

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

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

- Use the resulting bias-corrected parameter estimates to compute **(in)direct effects** for a change from **x̃ to x:**

$$OR^{NDE} \cong \frac{\exp\left(\theta_X x\right) \left\{1 + \exp\left(\theta_M + \theta_{XM} x + \beta_0 + \beta_X \tilde{x} + \beta_C c\right)\right\}}{\exp\left(\theta_X \tilde{x}\right) \left\{1 + \exp\left(\theta_M + \theta_{XM} \tilde{x} + \beta_0 + \beta_X \tilde{x} + \beta_C c\right)\right\}}$$

$$OR^{NIE} \cong \frac{\left\{1 + \exp\left(\beta_0 + \beta_X \tilde{x} + \beta_C c\right)\right\} 1 + \exp\left(\theta_X + \theta_{XM} x + \beta_0 + \beta_X x + \beta_C c\right)\right\}}{\left\{1 + \exp\left(\beta_0 + \beta_X x + \beta_C c\right)\right\} \left\{1 + \exp\left(\theta_M + \theta_{XM} x + \beta_0 + \beta_X \tilde{x} + \beta_C c\right)\right\}}$$

# We applied our methods to a preterm birth study

**?** 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**?
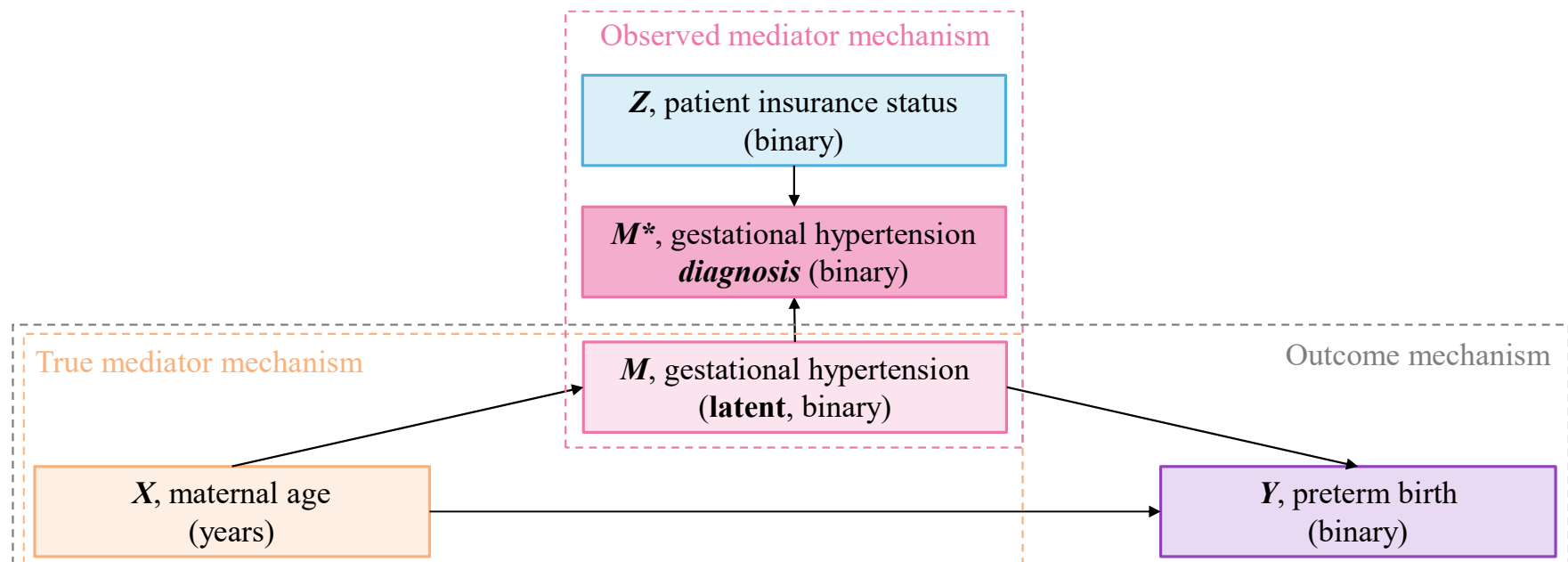
# We applied our methods to a preterm birth study

**Data:** National Vital Statistics System
- Provides demographic and health data for all births in a year in the US.
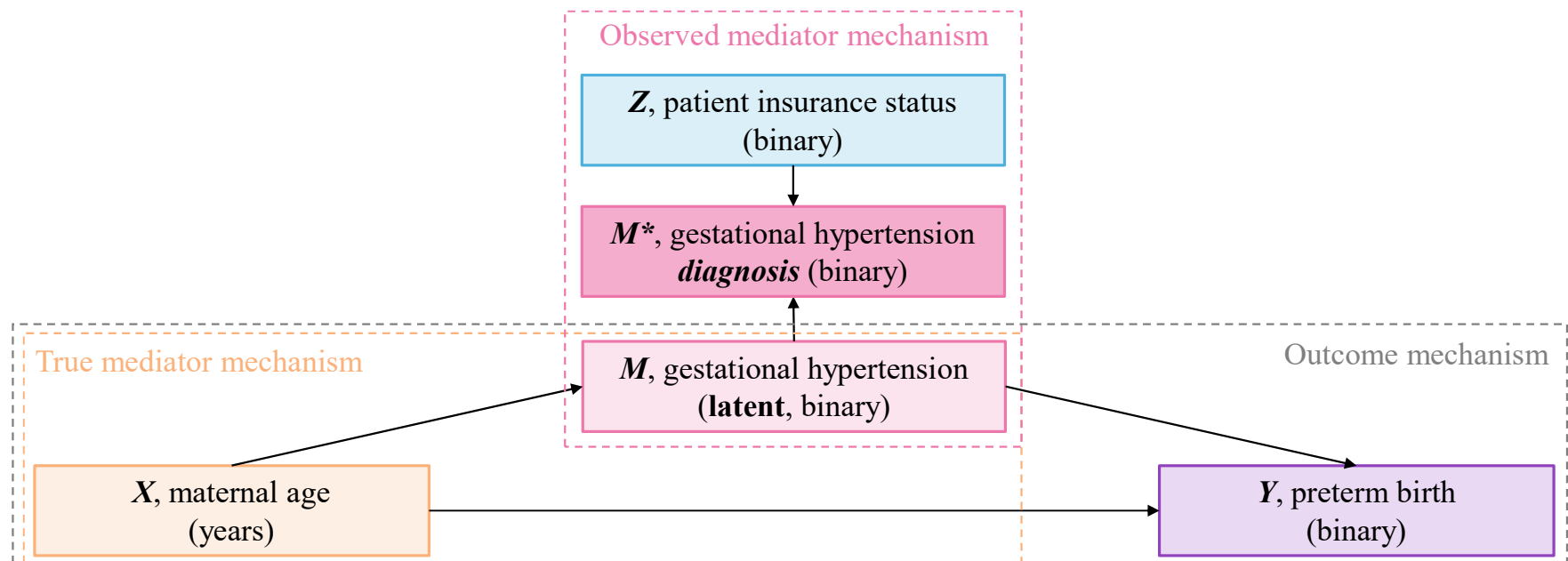- Random subsample from calendar year 2021, **N = 20,000**.

# We applied our methods to a preterm birth study

**True mediator mechanism:** $M \sim X$ + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** $M^* \mid M \sim$ Race + Z

**Outcome mechanism:** $Y \sim X$ + Race + Education + Parity + Smoking + BMI + $M$ + $M * X$

# Results change when we account for misdiagnosis

**True mediator mechanism:**  M ~ X + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** M* | M ~ Race + Z

**Outcome mechanism:** Y ~ X + Race + Education + Parity + Smoking + BMI + M + M * X

|  | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
|  | Est. | SE | Est. | SE |
| $\beta_X$ |  |  |  |  |
| $\gamma_{Z, G = 1}$ |  |  |  |  |
| $\gamma_{Z, G = 2}$ |  |  |  |  |
| $\theta_X$ |  |  |  |  |
| $\theta_M$ |  |  |  |  |
| $\theta_{XM}$ |  |  |  |  |

# Results change when we account for misdiagnosis

**True mediator mechanism:** $M$ ~ $X$ + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** $M^*$ | $M$ ~ Race + $Z$

**Outcome mechanism:** $Y$ ~ $X$ + Race + Education + Parity + Smoking + BMI + $M$ + $M * X$

Association between **age** & **gestational hypertension** is unchanged after accounting for misdiagnosis.

|  | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
|  | Est. | SE | Est. | SE |
| $\beta_X$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{Z, G = 1}$ |  |  |  |  |
| $\gamma_{Z, G = 2}$ |  |  |  |  |
| $\theta_X$ |  |  |  |  |
| $\theta_M$ |  |  |  |  |
| $\theta_{XM}$ |  |  |  |  |

# Results change when we account for misdiagnosis

**True mediator mechanism:** M ~ X + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** M* | M ~ Race + Z

**Outcome mechanism:** Y ~ X + Race + Education + Parity + Smoking + BMI + M + M * X

Association between **age** & **gestational hypertension** is unchanged after accounting for misdiagnosis.

Association between **gestational hypertension** & **preterm birth** strengthens.

|  | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
|  | Est. | SE | Est. | SE |
| $\beta_X$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{Z, G = 1}$ |  |  |  |  |
| $\gamma_{Z, G = 2}$ |  |  |  |  |
| $\theta_X$ | 0.02 | 0.05 | 0.10 | 0.03 |
| $\theta_M$ | 1.19 | 0.17 | 0.88 | 0.06 |
| $\theta_{XM}$ | 0.19 | 0.09 | 0.06 | 0.06 |

# Results change when we account for misdiagnosis

**True mediator mechanism:** M ~ X + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** M* | M ~ Race + Z

**Outcome mechanism:** Y ~ X + Race + Education + Parity + Smoking + BMI + M + M * X

| | EM Algorithm | | Naïve Analysis | |
|---|---|---|---|---|
| | Est. | SE | Est. | SE |
| $\beta_X$ | 0.10 | 0.04 | 0.08 | 0.03 |
| $\gamma_{Z, G=1}$ | -1.01 | 0.40 | - | - |
| $\gamma_{Z, G=2}$ | 2.09 | 8.81 | - | - |
| $\theta_X$ | 0.02 | 0.05 | 0.10 | 0.03 |
| $\theta_M$ | 1.19 | 0.17 | 0.88 | 0.06 |
| $\theta_{XM}$ | 0.19 | 0.09 | 0.06 | 0.06 |

Association between **age** & **gestational hypertension** is unchanged after accounting for misdiagnosis.

Association between **gestational hypertension** & **preterm birth** strengthens.

Use γ estimates to compute **sensitivity and specificity**.

45

# *M* is measured with perfect specificity and low sensitivity

**True mediator mechanism:** **M** ~ **X** + Race + Education + Parity + Smoking + BMI

**Observed mediator mechanism:** **M\* | M** ~ Race + **Z**

**Outcome mechanism:** **Y** ~ **X** + Race + Education + Parity + Smoking + BMI + **M** + **M** \* **X**

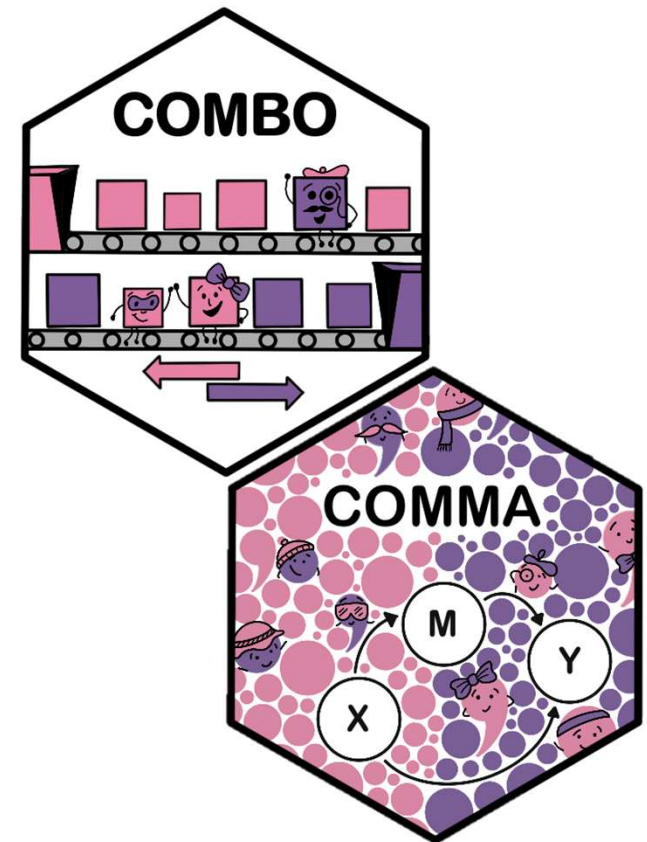|  | **Estimated Specificity**<br>P( no M\* \| no M ) | **Estimated Sensitivity**<br>P( M\* \| M ) |
|---|---|---|
| Insured | 99.9% | 43.1% |
| Self-Pay | 99.4% | 21.7% |

# 🗝 Key takeaways

- Developed new methods for handling misclassified binary mediator variables.
- Computed (in)direct effects using bias-corrected parameter estimates.
- Quantified gestational hypertension misdiagnosis rates based on insurance status.

# Software

- Estimation methods for **misclassified outcomes** are available in the *COMBO* R Package on CRAN.
  - **Co**rrecting **M**isclassified **B**inary **O**utcomes

- Estimation methods for **misclassified mediators** are available in the *COMMA* R Package on CRAN.
  - **Co**rrecting **M**isclassified **M**ediation **A**nalysis

# Thank you!

**Kimberly A. H. Webb**

kimberlywebb@pitt.edu

kimhwebb.com ⟶ My "webb-site" ☺

Scan for the paper

bit.ly/Webb-SMMR

University of **Pittsburgh**

Cornell Bowers CIS
**Statistics and Data Science**