# A **Primary Challenge** in Multi-site International Study

**Individual Patient-level Data (IPD) cannot be shared across sites**

- Regulatory Approval Processes
- Country-Specific Laws (e.g., HIPAA in USA, PIPEDA in Canada)
- Institutional Policies Data Sharing Restrictions

**Privacy-Preserving Federated Learning Algorithms**

- Enables fitting statistical models in a federated manner
- Requires summary statistics, instead of IPD
- Ensures data privacy and security

# Two Ideal Properties of Federated Learning Algorithms

**To date, only a few algorithms have successfully achieved both lossless and one-shot properties simultaneously:**

- Linear Regression (i.e., Chen et al., 2006, IEEE)
- Linear mixed models (i.e., Luo et al., 2022, Nature Communications)

Lossless

One-shot

Data Warehouse

$Y_1$

$Y_2$

$Y_3$

Results

Pool

Identical Results

No accuracy loss due to data sharing constraints

Only a single round of communication is needed

Luo C, Islam MN, Sheils NE, Buresh J, Reps J, Schuemie MJ, Ryan PB, Edmondson M, Duan R, Tong J, Marks-Anglin A. DLMM as a lossless one-shot algorithm for collaborative multi-site distributed linear mixed models. Nature communications. 2022 Mar 30;13(1):1678.

Yixin Chen et al. Regression Cubes with Lossless Compression and Aggregation. IEEE Trans Knowl Data Eng 18, 1585–1599 (2006).

September 16, 2025

# Challenges in Real-world Data



**Non-continuous outcomes**
- Binary outcome
- Categorical outcome
- Count outcome
- …

**Between-site heterogeneity**

We need **Federated Learning Algorithms** for **Generalized Linear Mixed Model (GLMM)**

**GLMM Model:** $\quad y_{ij} = b_i + \boldsymbol{\beta} x_{ij} + \epsilon_{ij}$, **where** $\epsilon_{ij} \sim N(0, \sigma^2), b_i \sim N(0, \tau^2)$ **is site-specific random effect for k-th site,** $\beta$ **is fixed effect.**

# Existing Works on Federated Learning Algorithms for GLMM

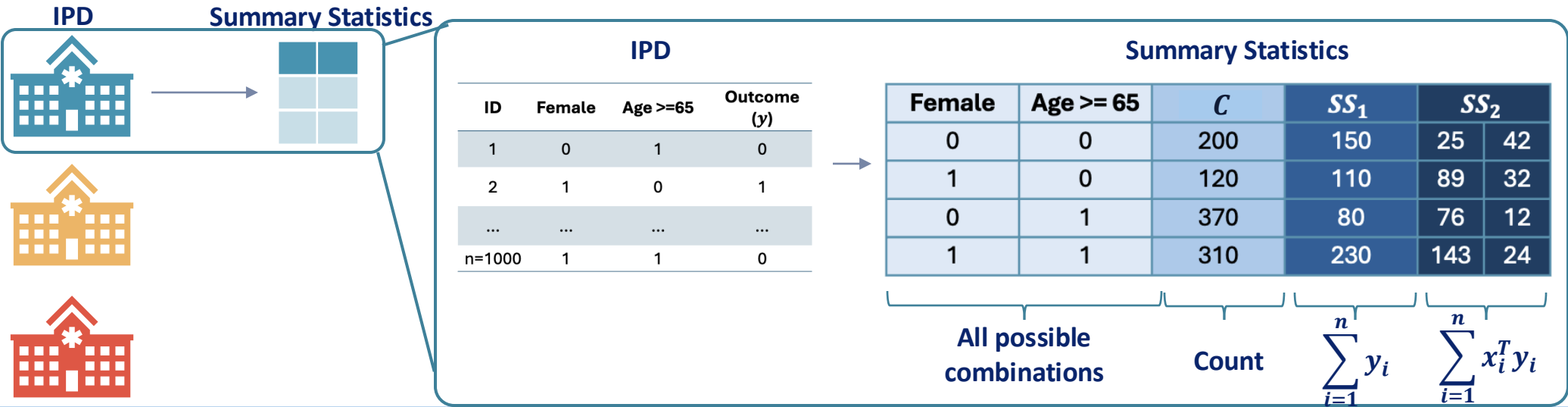| | Zhu et al, 2020, *Bioinformatics* | Luo et al, 2022, *JAMIA* | Yan et al, 2022, *arxiv* | COLA-GLMM |
|---|---|---|---|---|
| **Lossless** | ❌ | ✅ | ❌ | ✅ |
| **One-shot** | ❌ | ❌ | ⓘ | ✅ |
| **Communication Round** | Iterative (500~1000 rounds) | < 5 rounds | 1 or 2 rounds* (Depends on initialization) | 1 round |

# Proposed Method – COLA-GLMM

**Collaborative One-shot Lossless Algorithm for Generalized Linear Mixed Model**
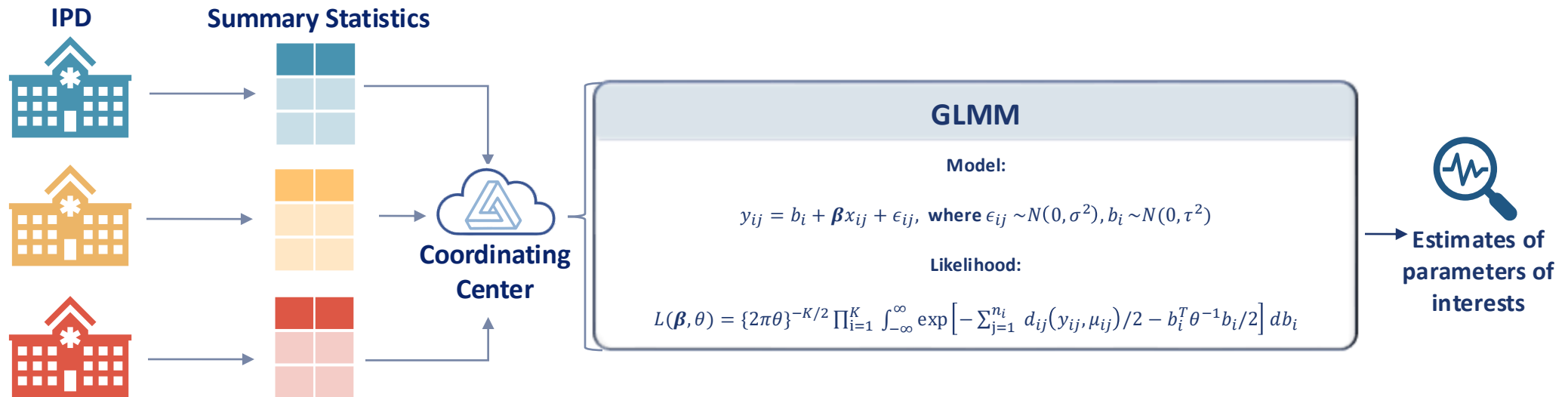
- Suppose that a **common** set of covariates are available at all collaborating sites.

- The covariates have been **standardized into categorical variables**.

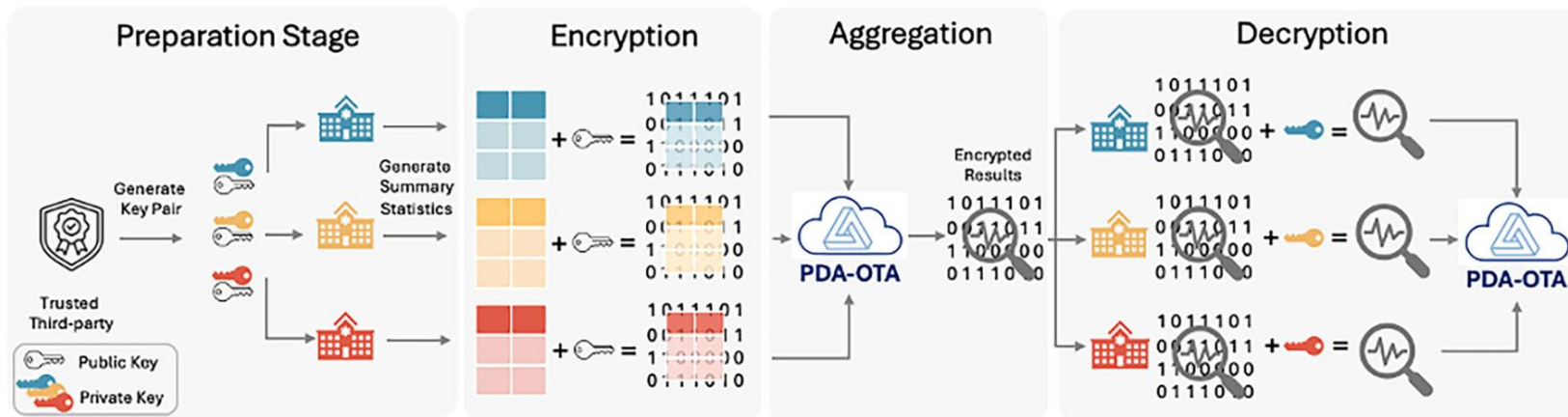- Pipeline:

# Proposed Method – COLA-GLMM

**Collaborative One-shot Lossless Algorithm for Generalized Linear Mixed Model**

- Suppose that a **common** set of covariates are available at all collaborating sites.

- The covariates have been **standardized into categorical variables**.

- Pipeline:



**IPD**   **Summary Statistics**

**Coordinating Center**

**GLMM**

**Model:**

$$y_{ij} = b_i + \boldsymbol{\beta} x_{ij} + \epsilon_{ij}, \text{ where } \epsilon_{ij} \sim N(0, \sigma^2), b_i \sim N(0, \tau^2)$$

**Likelihood:**

$$L(\boldsymbol{\beta}, \theta) = \{2\pi\theta\}^{-K/2} \prod_{i=1}^{K} \int_{-\infty}^{\infty} \exp\left[ -\sum_{j=1}^{n_i} d_{ij}(y_{ij}, \mu_{ij})/2 - b_i^T \theta^{-1} b_i / 2 \right] db_i$$

**Estimates of parameters of interests**

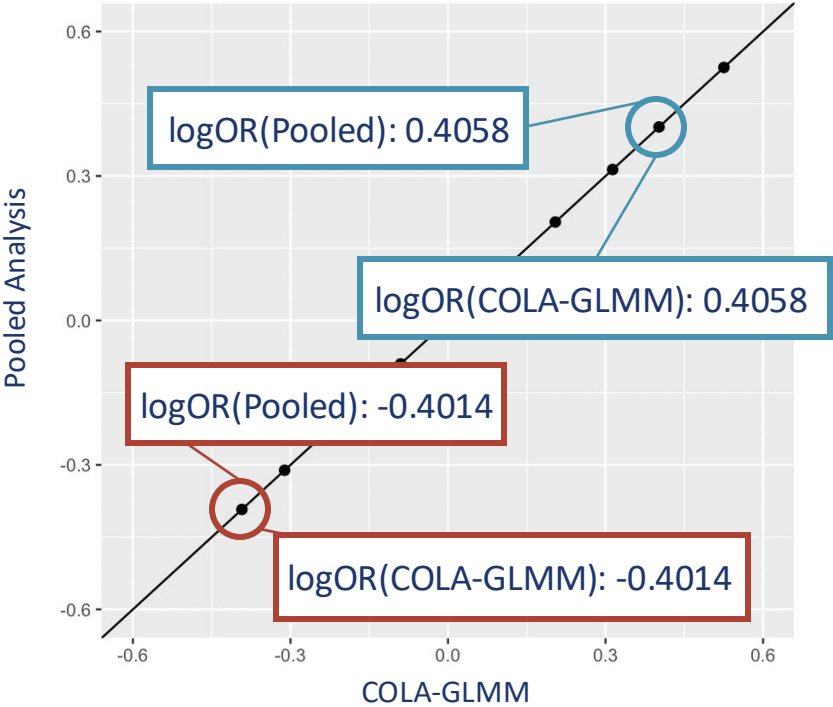# Homomorphic Encryption Enhanced COLA-GLMM

- Under semi-trusted environment:
  - while both data contributors and the coordinating center adhere to the protocol without engaging in malicious actions, the coordinating center may still attempt to derive insights from passively obtained data, indicating a curiosity in extracting information.

# Simulation Study – Compare Pooled Analysis and COLA-GLMM

## No cell suppression



logOR(Pooled): 0.4058

logOR(COLA-GLMM): 0.4058

logOR(Pooled): -0.4014

logOR(COLA-GLMM): -0.4014

Pooled Analysis

COLA-GLMM

---

U.S. Dept. of Health & Human Services
**Guidance Portal**

Return to Search

## CMS Cell Suppression Policy

Guidance for CMS Cell Suppression Policy Web Page

Final

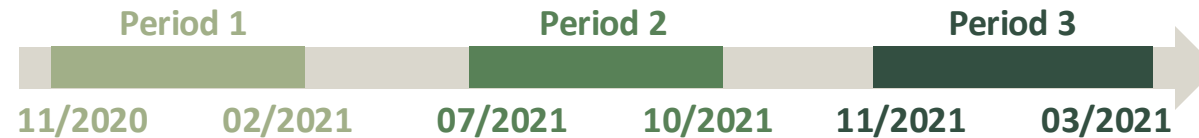**Issued by:** Centers for Medicare & Medicaid Services (CMS)

| Raw value | | Report |
|-----------|---|--------|
| 0 | | 0 |
| 1-10 | | <11 |
| >=11 | | as it |

# Real-world Case Study

- **Scientific Question:**

  Identify COVID-19 mortality **risk factors**
  over **three time periods** among hospitalized patients

- **Study Period:**

  Period 1      Period 2      Period 3

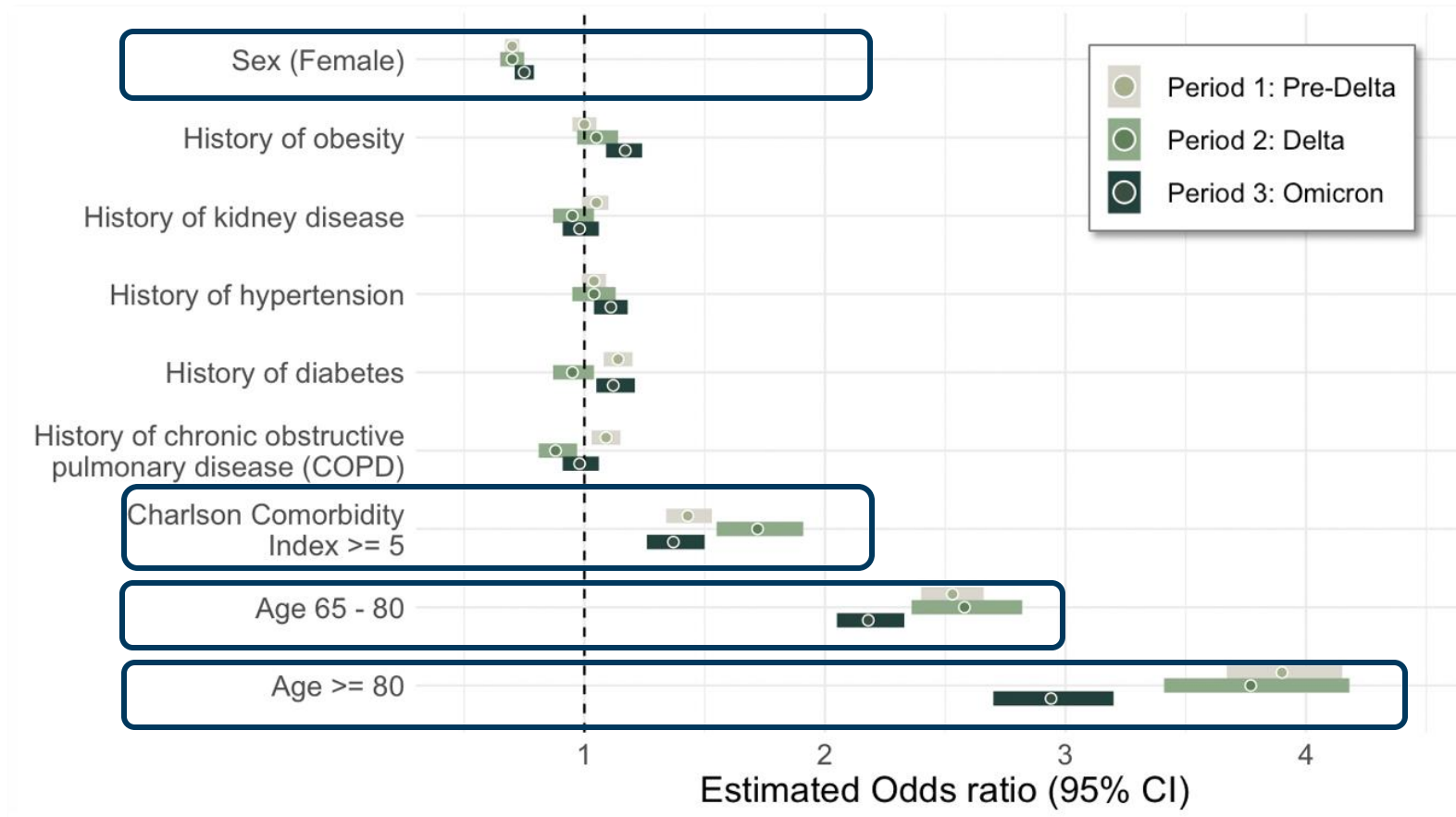  **11/2020**   **02/2021**   **07/2021**   **10/2021**   **11/2021**   **03/2021**

- **Databases (3 countries):**

  - Optum® de-identified Electronic Health Record Dataset (Optum EHR);
  - Optum's Clinformatics® Data Mart (CDM or Clinformatics®);
  - IQVIA Hospital CDM;
  - University of Florida Health;
  - Department of Veterans Affairs;
  - Integrated Primary Care Information (IPCI), The Netherlands;
  - Columbia University Irving Medical Center (CUIMC);
  - Parc Salut Mar Barcelona (PSMAR), Spain.

**Inclusion criteria:**

- Patients aged 18 years and older
- Had an inpatient visit with either a diagnosis of COVID-19 or a positive test for COVID-19 between 21 days prior to the inpatient visit and the end of the inpatient visit

September 16, 2025

# Real-world Case Study Results



- **Sex (female):**
  - Reference group: Male
  - Female patients consistently exhibit a lower risk of mortality compared to males across all periods

- **Charlson Comorbidity Index (CCI):**
  - Reference group: CCI < 5
  - Higher CCI scores are statistically associated with an increased risk of mortality.

- **Age:**
  - Reference group: Age < 65
  - Higher age indicates significantly increased risk of mortality

# Summary – COLA-GLMM

**C**ollaborative **O**ne-shot **L**ossless **A**lgorithm for **G**eneralized **L**inear **M**ixed **M**odel

- **Lossless One-Shot**
- **Summary Statistics Only**
- **Heterogeneity-Aware**
- **Scalable, Applicable, and Implementation-Ready in OHDSI Network**



**PDA R Package: 13300+ downloads since 2020**

**PDA Github Page: https://github.com/Penncil/pda**

**PDA website: https://pdamethods.org/**

**PDA-OTA: https://pda-ota.pdamethods.org/**
Penn security office certified

# Acknowledgements

September 16, 2025