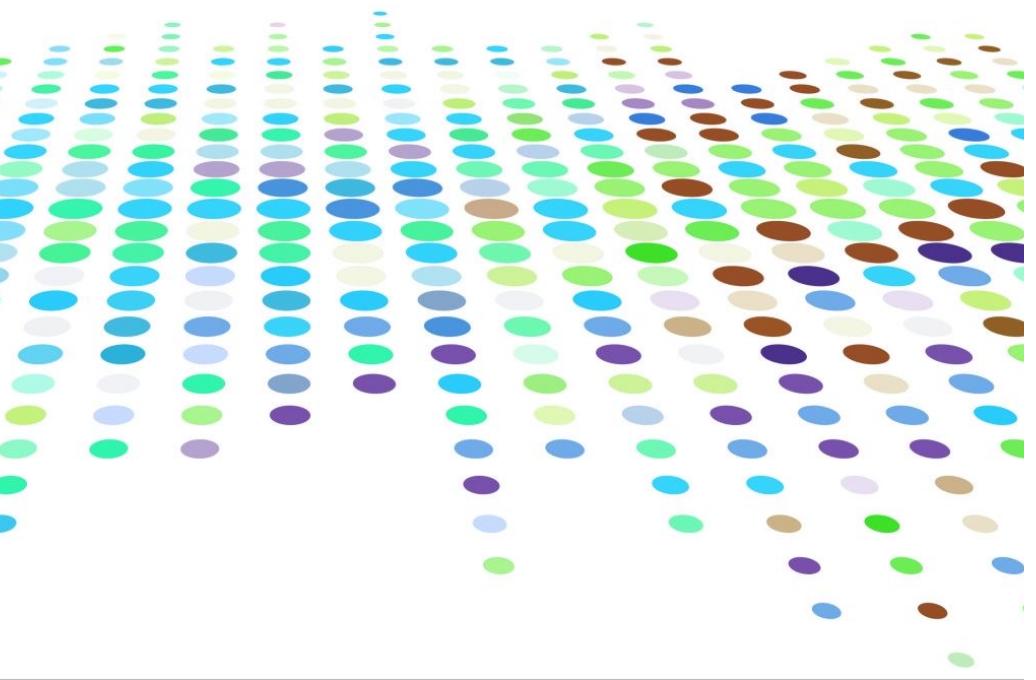

PRINCIPLES AND PRACTICES OF FEDERATED LEARNING: METHODS, CHALLENGES, AND CASE STUDIES



Rui Duan (段芮)

Harvard T.H. Chan School of Public Health

2024 International Conference on Frontiers of Data
Science

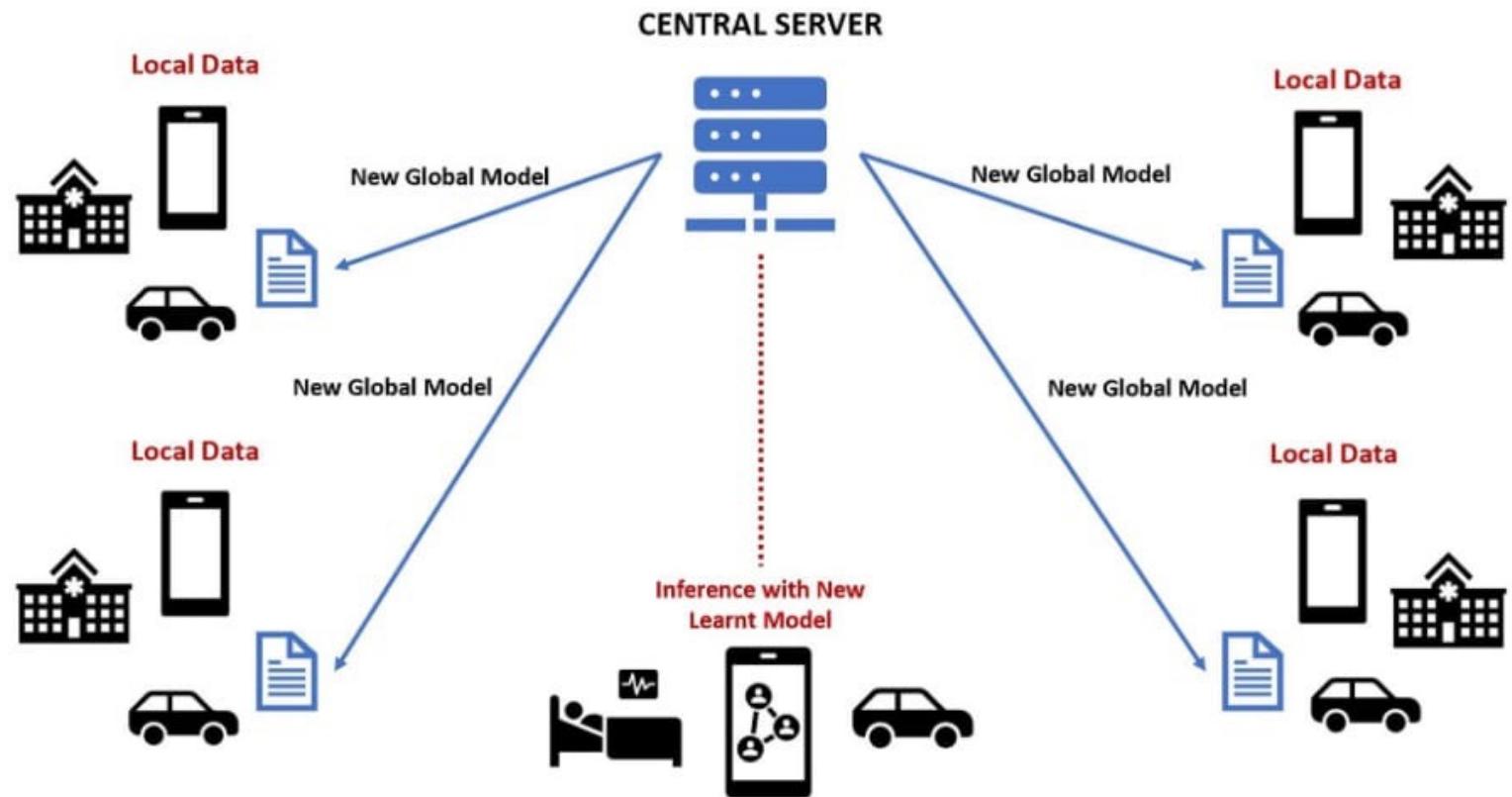
OUTLINE

- Section I — Overview of Federated Learning
 - What is Federated Learning?
 - Important applications
 - Challenges
 - Section II — Federated Learning via Federated Optimization
 - Section III — Federated Statistical Inference
 - Discussion— Other Related Topics
-

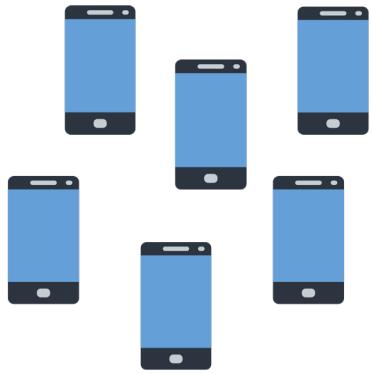
SECTION I: OVERVIEW

FEDERATED LEARNING

- Federated Learning is a machine learning technique that allows training models across multiple decentralized devices or servers holding local data samples without exchanging them.



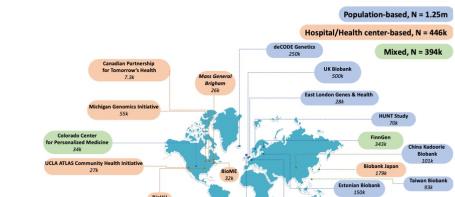
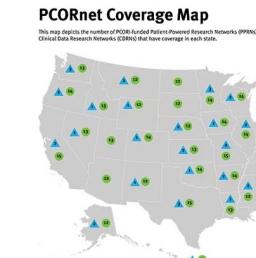
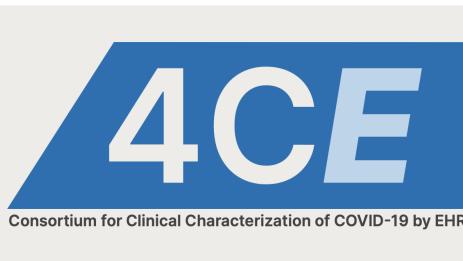
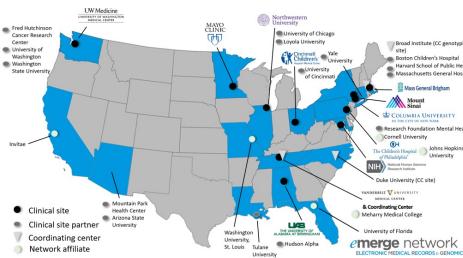
Picture from internet



Datafication: Transforming the World into Data

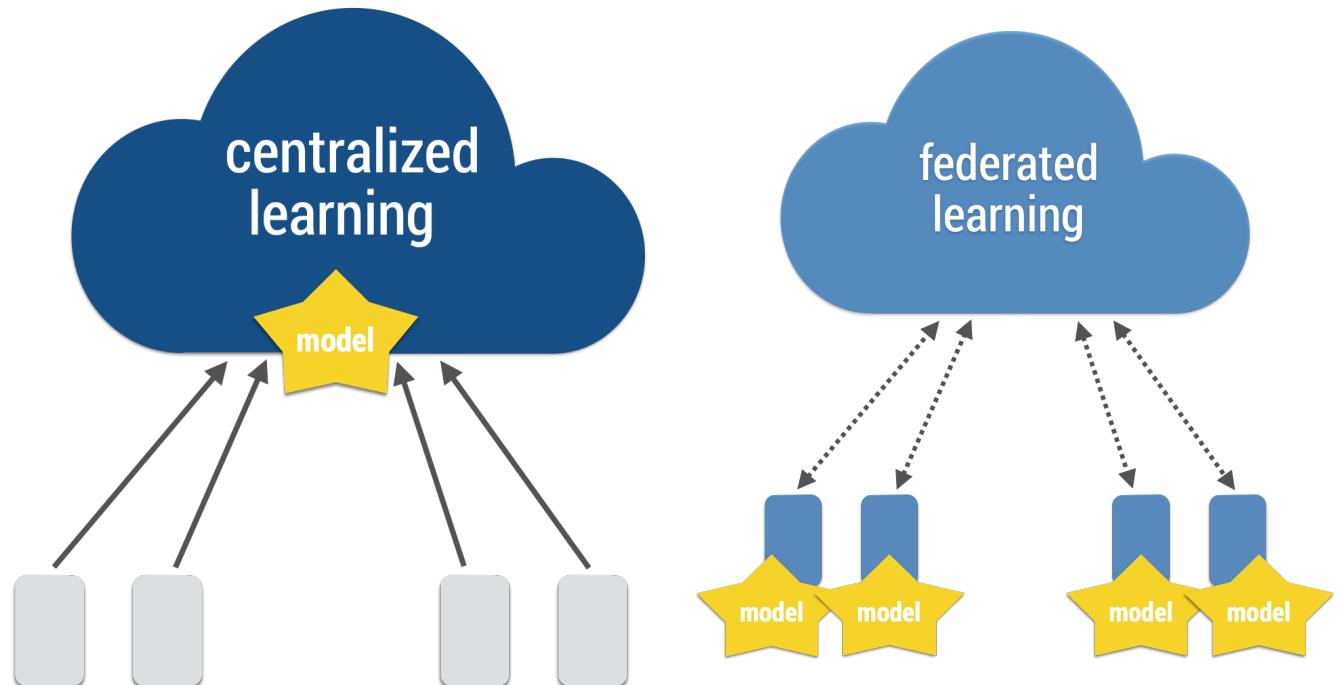
Pictures from internet

DATA NETWORKS

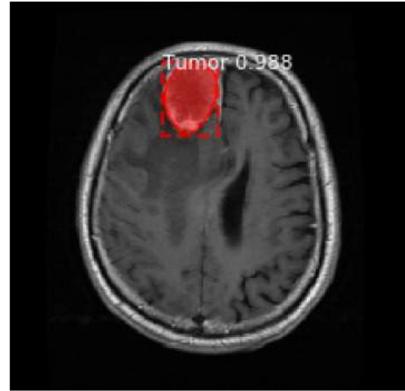


WHY FEDERATED LEARNING

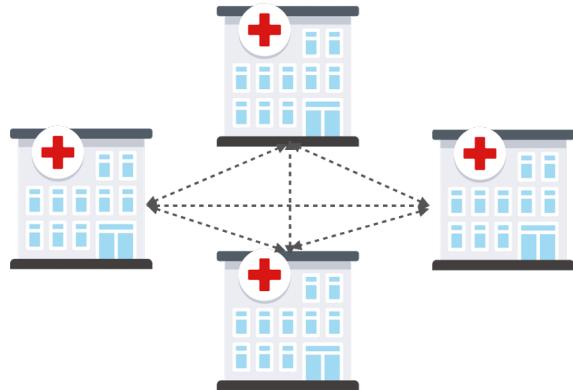
- *reduce strain on network*
- *incorporate more data*
- *improve privacy**



Picture from Federated and Collaborative Learning, Fall 2023, Virginia Smith.
<https://www.cs.cmu.edu/~smithv/10719/>



collaboratively learn tumor detection
model without sharing patient data



NVIDIA.

PLATFORMS ▾ DEVELOPERS ▾ INDUSTRIES ▾ SHOP DRIVERS ▾ SUPPORT ABOUT NVIDIA ▾ EMAIL SIGN-UP

HOME DEEP LEARNING NETWORKING DRIVING GAMING PRO GRAPHICS AUTONOMOUS MACHINES HEALTHCARE AI PODCAST

Medical Institutions Collaborate to Improve Mammogram Assessment AI with NVIDIA Clara Federated Learning

In a federated learning collaboration, the American College of Radiology, Diagnostics da America, Partners HealthCare, Ohio State University and Stanford Medicine developed better predictive models to assess breast tissue density.

April 15, 2020 by MONA FLORES

ARTIFICIAL INTELLIGENCE, DIAGNOSTICS

UPenn, Intel partner to use federated learning AI for early brain tumor detection

The project will bring in 29 institutions from North America, Europe and India and will use privacy-preserved data to train AI models. Federated learning has been described as being born at the intersection of AI, blockchain, edge computing and the Internet of Things.

By ALARIC DEARMANT

Post a comment / May 11, 2020 at 10:03 AM

LILY HAY NEWMAN SECURITY 10.07.2020 02:19 PM

How Google's Android Keyboard Keeps 'Smart Replies' Private

The latest Gboard feature needs to know as much as possible about your digital life to work—but doesn't share that data with Google.

Artificial intelligence / Machine learning

How Apple personalizes Siri without hoovering up your data

The tech giant is using privacy-preserving machine learning to improve its voice assistant while keeping your data on your phone.

by Karen Hao December 11, 2019

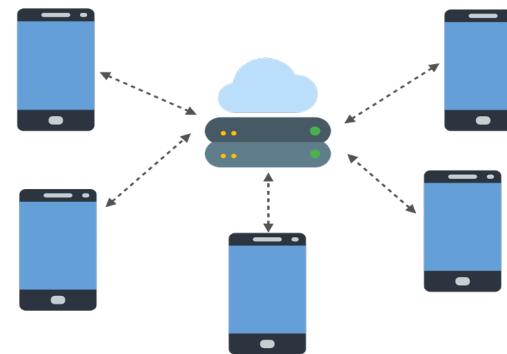
Sherpa raises \$8.5M to expand from conversational AI to B2B privacy-first federated learning services

Ingrid Lunden @ingridlunden 7:11 PM EDT • March 15, 2021

Comment

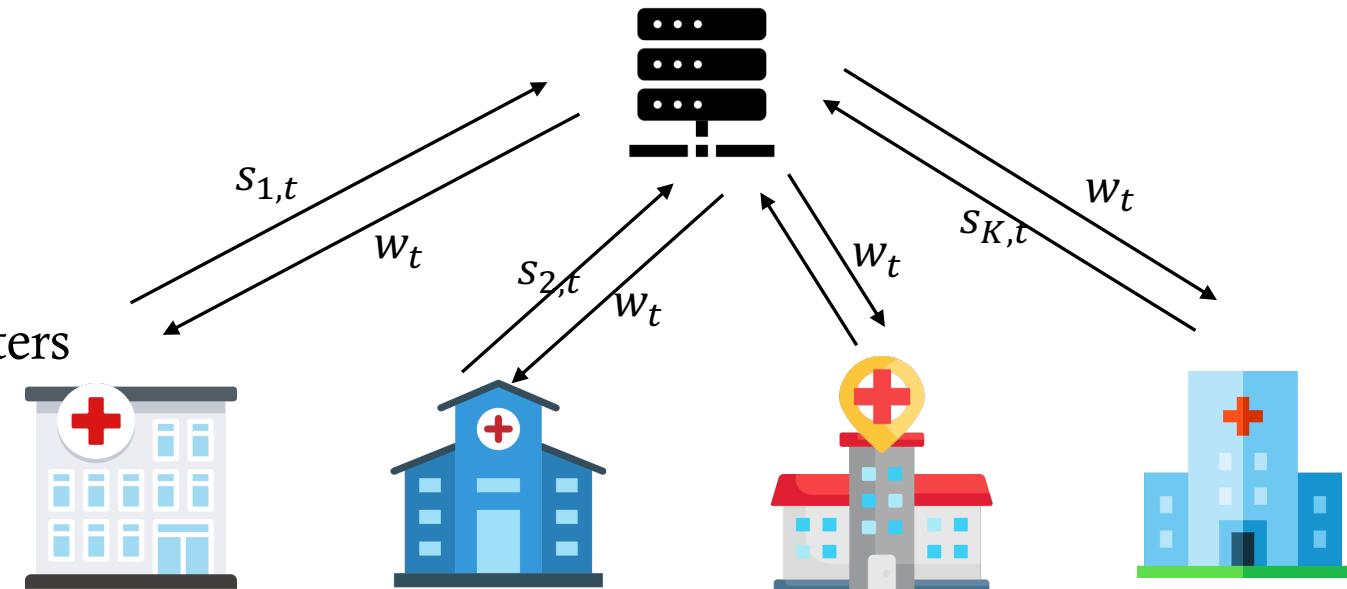


power next-word prediction in mobile keyboards without collecting texts

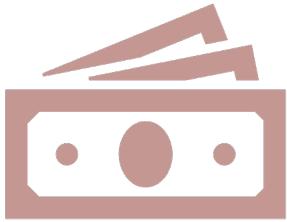


BASIC PROBLEM SET-UP

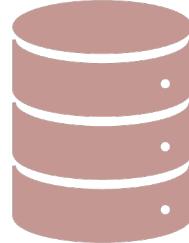
- Objective function: $L(w) = \sum_{k=1}^K L_k(w; D_k)$
- Optimization over all sites: $\hat{w} = \operatorname{argmin}_w L(w)$
- Initialization
- Local sites share summary stats
with central server
- Central server updates model parameters



CHALLENGES



Expensive
Communications



Data heterogeneity



Privacy concerns

COMMUNICATION COST



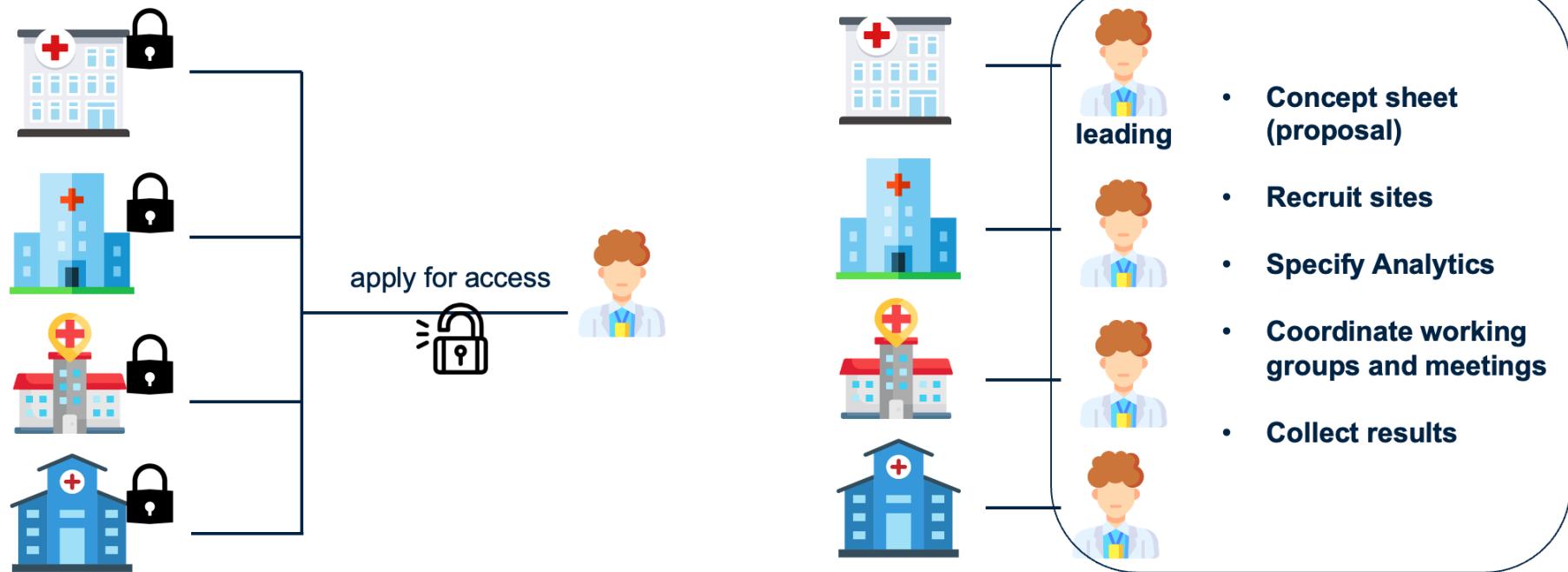
Upload and Download Costs

Limit on file sizes
Time

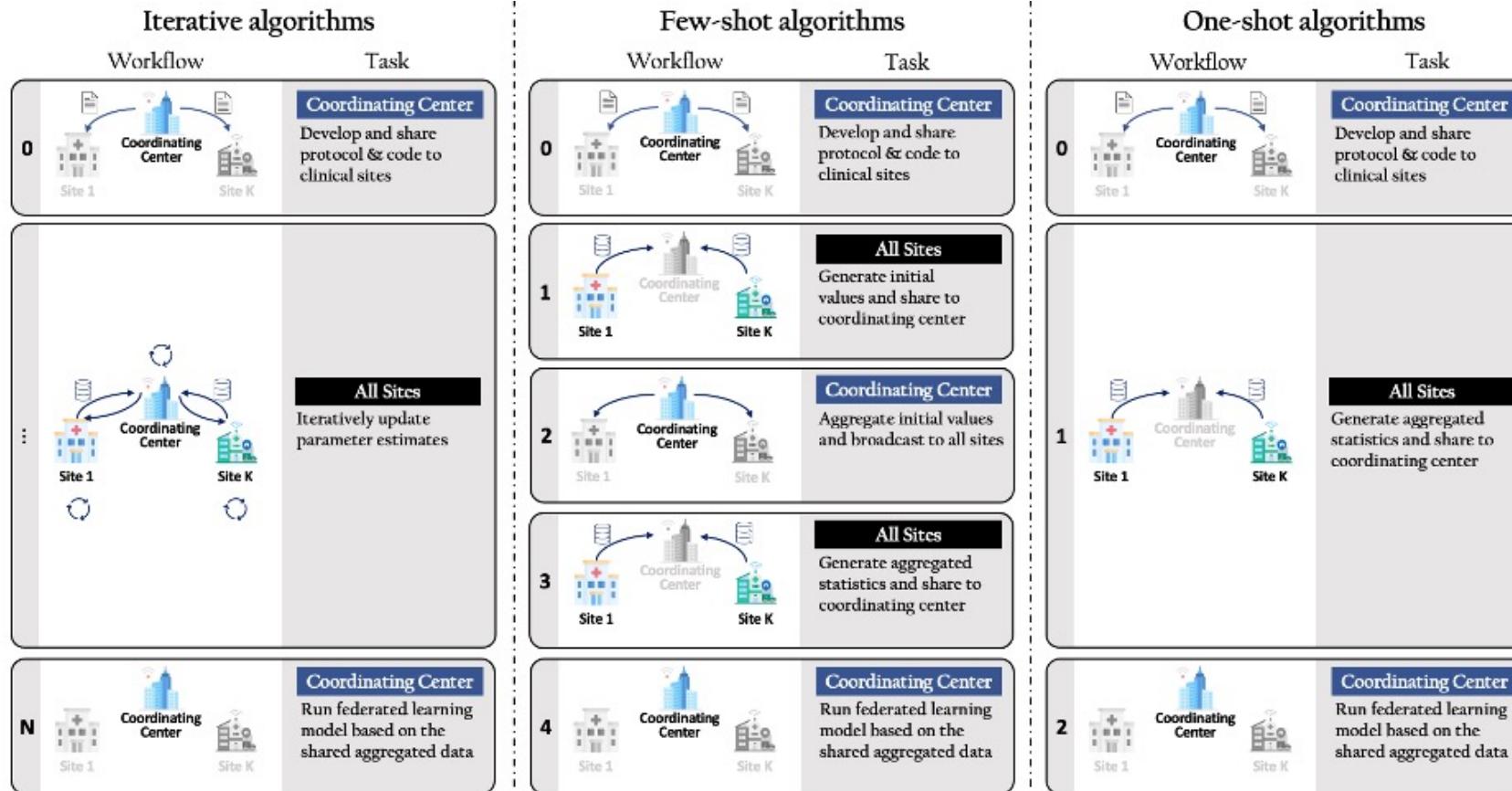


Human efforts involved in the communication

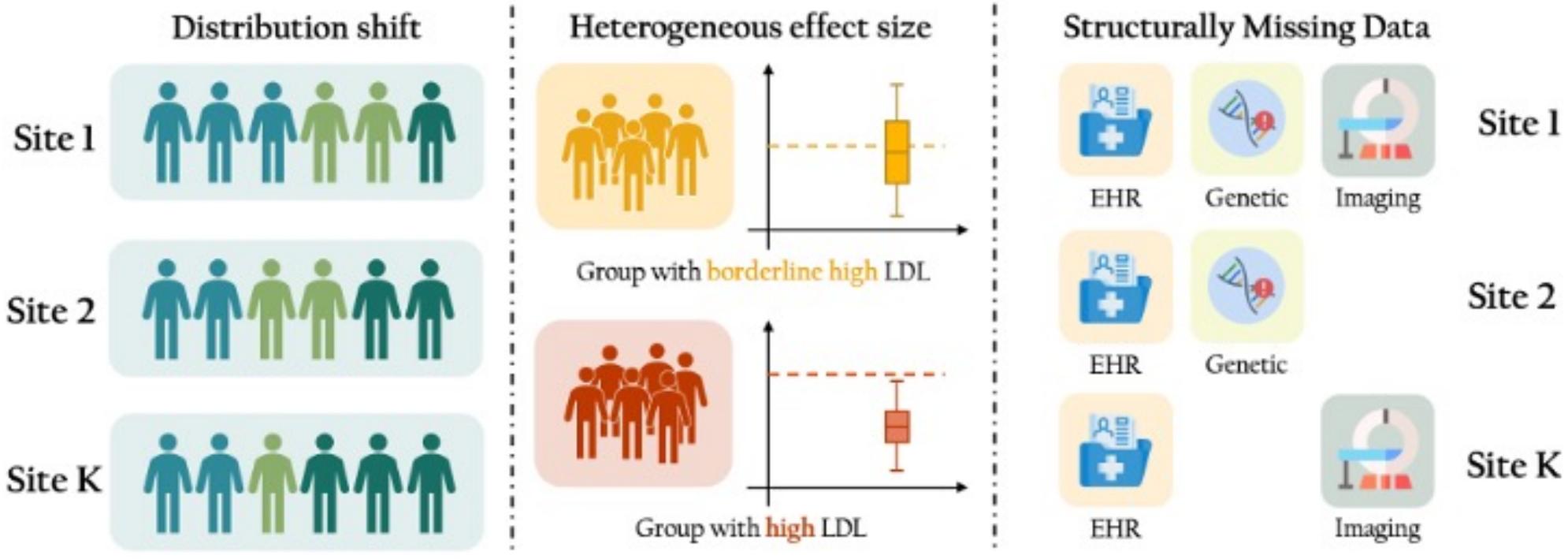
DATA SHARING IN RESEARCH NETWORK



REDUCING THE COMMUNICATION COST



DATA HETEROGENEITY

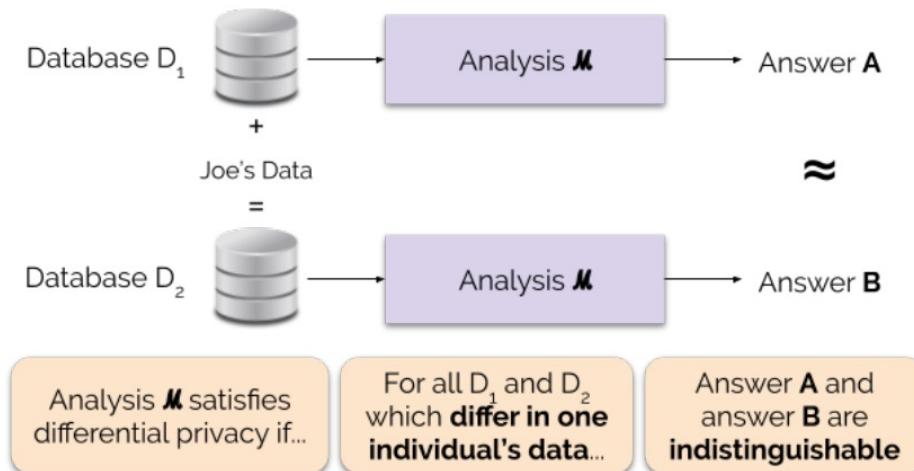


VARIOUS TASKS IN FEDERATED LEARNING

- Train a unified model for all sites
 - Better generalizability and transferability
- Jointly train site-specific models
 - Leveraging between-site similarity while accounting for site-level differences.
- Train a model for a target population of interest
 - Target population can be defined within or across sites
 - Leveraging all data within the network

PRIVACY CONCERNS

- Encryption
- Secure multiparty computation
- Differential privacy



intuition: can't determine whether or not Joe was present in the dataset,
let alone the contents of his data

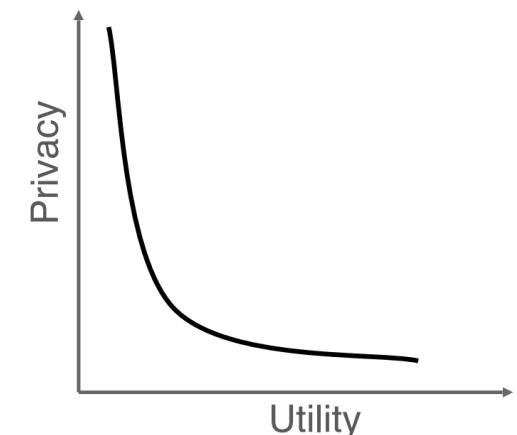
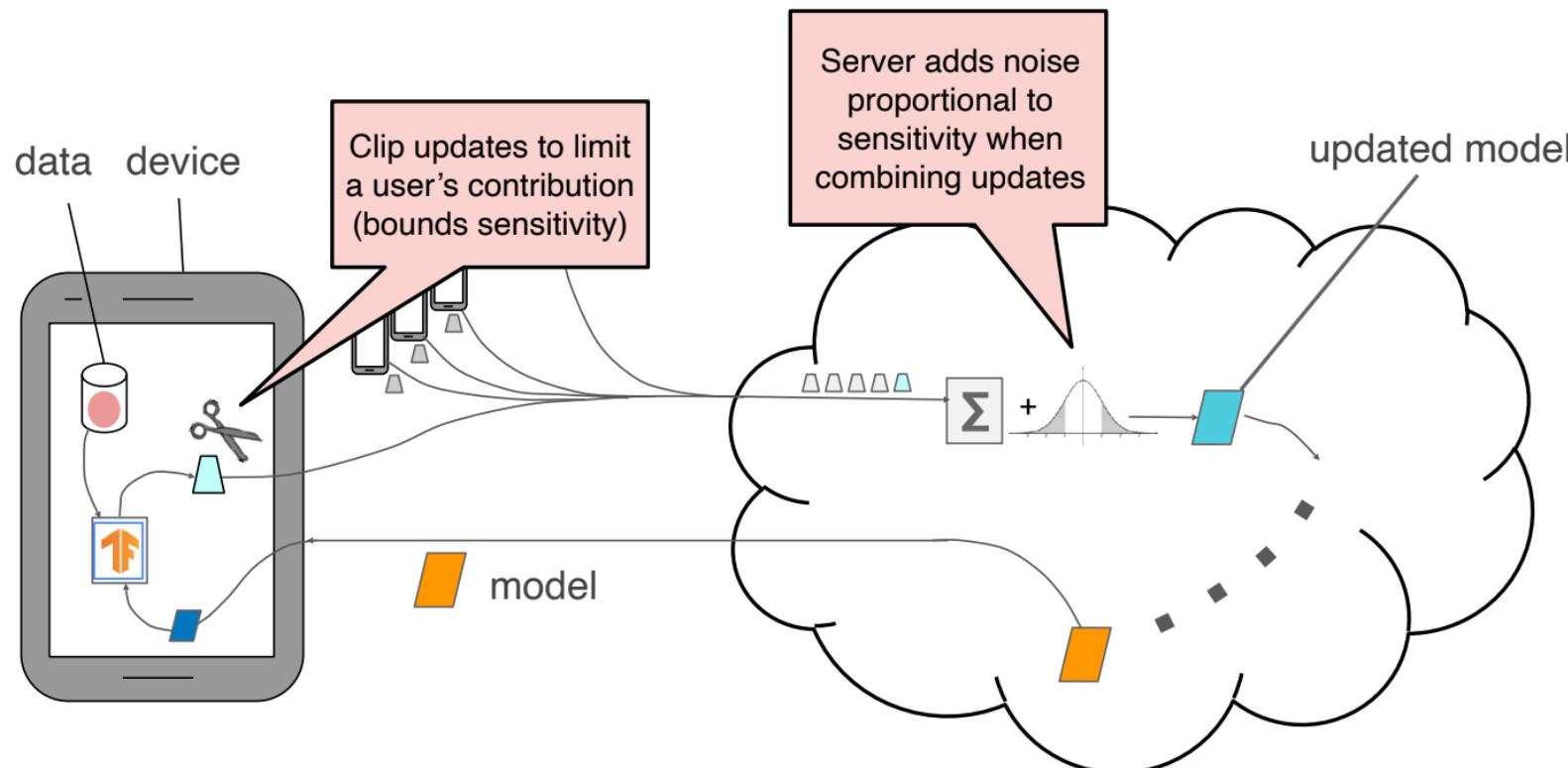
Picture from Federated and Collaborative Learning, Fall 2023, Virginia Smith.
<https://www.cs.cmu.edu/~smithv/10719/>

Ma, Jing, et al. "Privacy-preserving federated learning based on multi-key homomorphic encryption." *International Journal of Intelligent Systems* 37.9 (2022): 5880-5901.

Zhang, Chen, et al. "A survey on federated learning." *Knowledge-Based Systems* 216 (2021): 106775.

Li, Yong, et al. "Privacy-preserving federated learning framework based on chained secure multiparty computing." *IEEE Internet of Things Journal* 8.8 (2020): 6178-6186.

PRIVACY CONCERN



[Credit: P. Kairouz, B. McMahan, V. Smith, FL Tutorial, NeurIPS 2020]

A General Approach to Adding Differential Privacy to Iterative Training Procedures

H. Brendan McMahan
mcmahan@google.com

Galen Andrew
galenandrew@google.com

Úlfar Erlingsson
ulfar@google.com

Steve Chien
schien@google.com

Ilya Mironov
mironov@google.com

Nicolas Papernot
papernot@google.com

Peter Kairouz
kairouz@google.com

Abstract

In this work we address the practical challenges of training machine learning models on privacy-sensitive datasets by introducing a modular approach that minimizes changes to training algorithms, provides a variety of configuration strategies for the privacy mechanism, and then isolates and simplifies the critical logic that computes the final privacy guarantees. A key challenge is that training algorithms often require estimating many different quantities (vectors) from the same set of examples — for example, gradients of different layers in a deep learning architecture, as well as metrics and batch normalization parameters. Each of these

Introducing TensorFlow Privacy: Learning with Differential Privacy for Training Data



TensorFlow

Follow

Mar 6 · 7 min read



Posted by [Carey Radebaugh](#) (Product Manager) and [Úlfar Erlingsson](#) (Research Scientist)

Today, we're excited to announce TensorFlow Privacy ([GitHub](#)), an open source library that makes it easier not only for developers to train machine-learning models with privacy, but also for researchers to advance the state of the art in machine learning with strong privacy guarantees.

Modern machine learning is increasingly applied to create amazing new technologies and user experiences, many of which involve training

[Credit: P. Kairouz, B. McMahan, V. Smith, FL Tutorial, NeurIPS 2020]

BIGGER PICTURE: COLLABORATIVE LEARNING

- *learning across multiple sources, stakeholders*
- ***how to incentivize participation***
- ***develop trustworthy learning schemes***
- ***effectively leveraging heterogeneous sources***

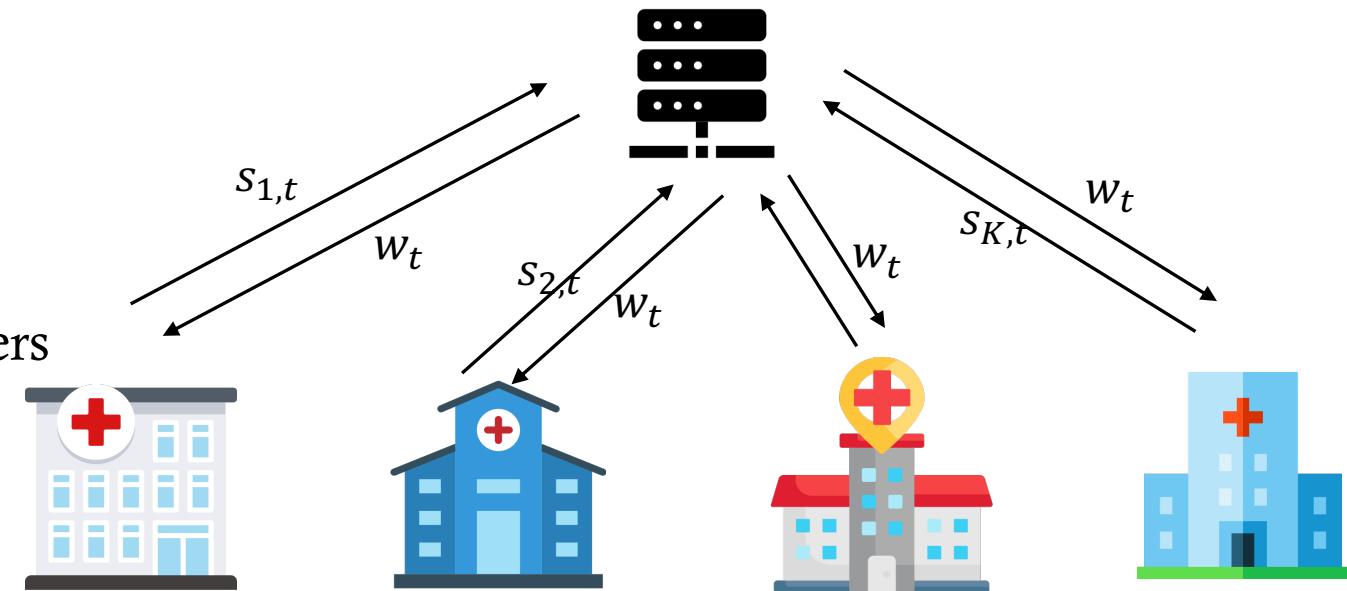


Pictures from internet

SECTION II: FEDERATED LEARNING VIA FEDERATED OPTIMIZATION

BASIC PROBLEM SET-UP

- Objective function: $L(w) = \sum_{k=1}^K L_k(w; D_k)$
- Optimization over all sites: $\hat{w} = \operatorname{argmin}_w L(w)$
- Initialization
- Local sites share summary stats
with central server
- Central server update model parameters



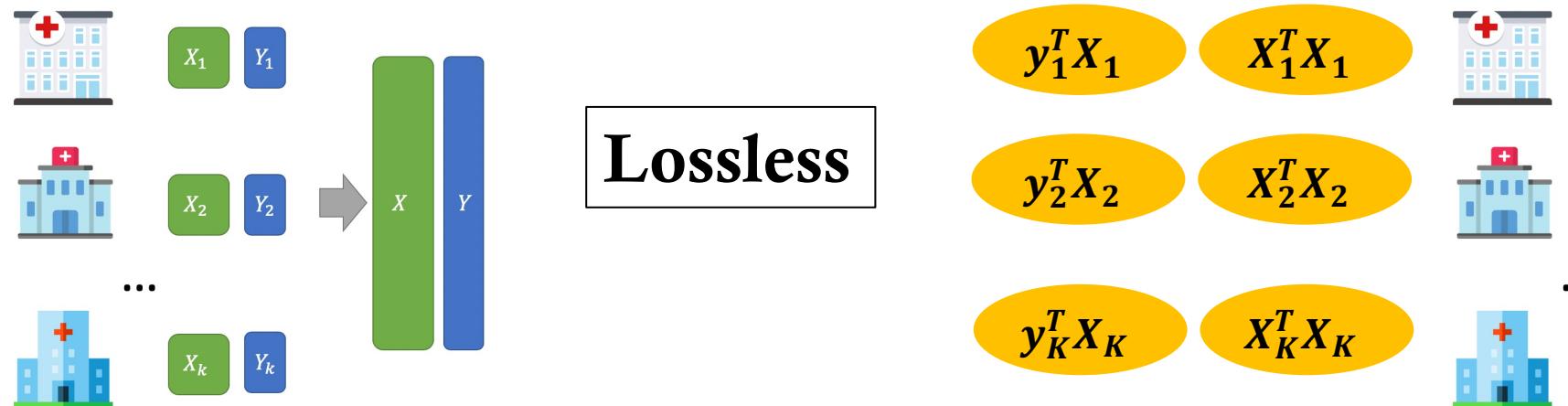
HOWTO REDUCE COMMUNICATION COST

- Leveraging properties of the loss function
 - Good initialization
 - Local updates
 - Trade-offs between information shared per communication and the number of communications
-

ONE-SHOT FEDERATED LEARNING

- Linear model

$$L(\beta) = \sum_{k=1}^K \|y_k - X_k\beta\|^2 = y_k^T y_k + \beta^T X_k^T X_k \beta - 2\beta^T X_k y_k$$



Chen et al.(2006) Regression cubes with lossless compression and aggregation.

LINEAR MIXED MODEL

- y_{ki} is the continuous outcome, $k = 1, \dots, K, i = 1, \dots, n_k$,
- x_{ki} is the p -dim covariates vector with β the fixed effects,
- z_{ki} is the q -dim covariates vector with u_k the random effects. z_{ki} usually is a subset of x_{ki} .
- $y_k = (y_{k1}, \dots, y_{kn_k})^T, X_k = (x_{k1}, \dots, x_{kn_k})^T, Z_k = (z_{k1}, \dots, z_{kn_k})^T,$
$$y_k = X_k\beta + Z_k u_k + \epsilon_k,$$
with $u_k \sim N(0, V), V = diag(\sigma_1^2, \dots, \sigma_q^2)$, random error $\epsilon_k \sim N(0, \sigma^2 I_{n_k})$.
- Need to estimate β , and variance components $\theta = (\sigma^2, \sigma_1^2, \dots, \sigma_q^2)$.

-
- Log-likelihood,

$$\ell(\beta, \theta) = -\frac{1}{2} \sum_{k=1}^K \{ \log |\Sigma_k| + (y_k - X_k \beta)^T \Sigma_k^{-1} (y_k - X_k \beta) \},$$

where $\Sigma_k = \Sigma_k(\theta) = Z_k V Z_k^T + \sigma^2 I_{n_k}$.

DLMM

- LMMs don't have closed-form solution
- Woodbury matrix identity

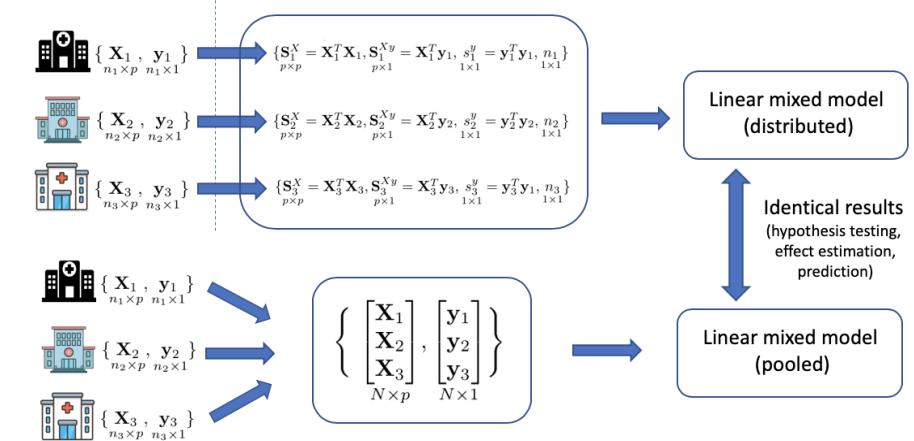
$$\Sigma_k^{-1} = \sigma^{-2} \{ I_{n_k} - Z_k (\sigma^2 V^{-1} + Z_k^T Z_k)^{-1} Z_k^T \}$$

- Matrix determinant lemma

$$|\Sigma_k| = \sigma^{2n_k} |I_q + \sigma^{-2} V Z_k^T Z_k|$$

- Log-likelihood,

$$\ell(\beta, \theta) = -\frac{1}{2} \sum_{k=1}^K \{ \log |\Sigma_k| + (y_k - X_k \beta)^T \Sigma_k^{-1} (y_k - X_k \beta) \}$$



nature communications

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [nature communications](#) > [articles](#) > [article](#)

Article | [Open Access](#) | Published: 30 March 2022

DLMM as a lossless one-shot algorithm for collaborative multi-site distributed linear mixed models

ITERATIVE ALGORITHM: GOOD INITIALIZATION

Initialize from
local site

Leveraging pre-
trained models

Initialize by
averaging across
all local models

Nguyen, John, et al. "Where to begin? on the impact of pre-training and initialization in federated learning." *arXiv preprint arXiv:2206.15387* (2022).

ITERATIVE ALGORITHM: FEDERATED AVERAGING (FEDAVG) VERSUS DISTRIBUTED GD (SGD)

- FedAvg: Locally updating models and taking the average of all local updates to get the global update
- Distributed GD: Sequentially update model across sites.

Distributed SGD: computation on device k

```
for i ∈ mini-batch B
| Δw ← Δw - α∇fi(w)
end
w ← w + Δw
```

FedAvg: computation on device k

```
for t = 1, 2, ..., local iterations T
| Δw ← Δw - α∇fit(w)
| w ← w + Δw
end
```

Why is it useful to perform `local-updating`?

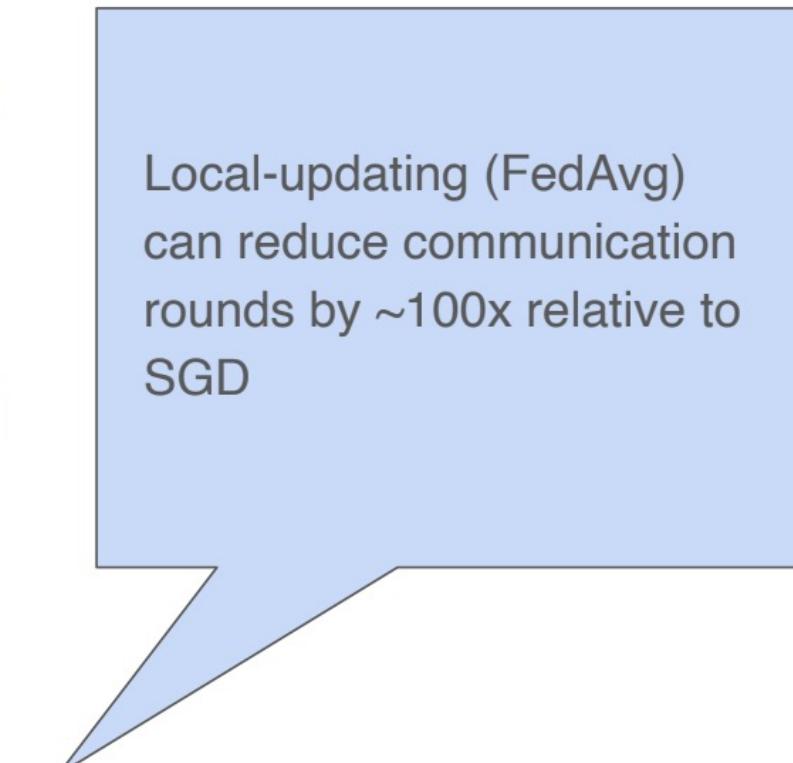
1. Can perform more local computation (i.e., more than just one mini-batch)
 2. Incorporate updates more quickly (immediately apply gradient information)
- ✓ Can lead to method converging in many fewer communication rounds

MNIST CNN, 99% ACCURACY

CNN	E	B	u	IID	NON-IID
FEDSGD	1	∞	1	626	483
FEDAVG	5	∞	5	179 (3.5x)	1000 (0.5x)
FEDAVG	1	50	12	65 (9.6x)	600 (0.8x)
FEDAVG	20	∞	20	234 (2.7x)	672 (0.7x)
FEDAVG	1	10	60	34 (18.4x)	350 (1.4x)
FEDAVG	5	50	60	29 (21.6x)	334 (1.4x)
FEDAVG	20	50	240	32 (19.6x)	426 (1.1x)
FEDAVG	5	10	300	20 (31.3x)	229 (2.1x)
FEDAVG	20	10	1200	18 (34.8x)	173 (2.8x)

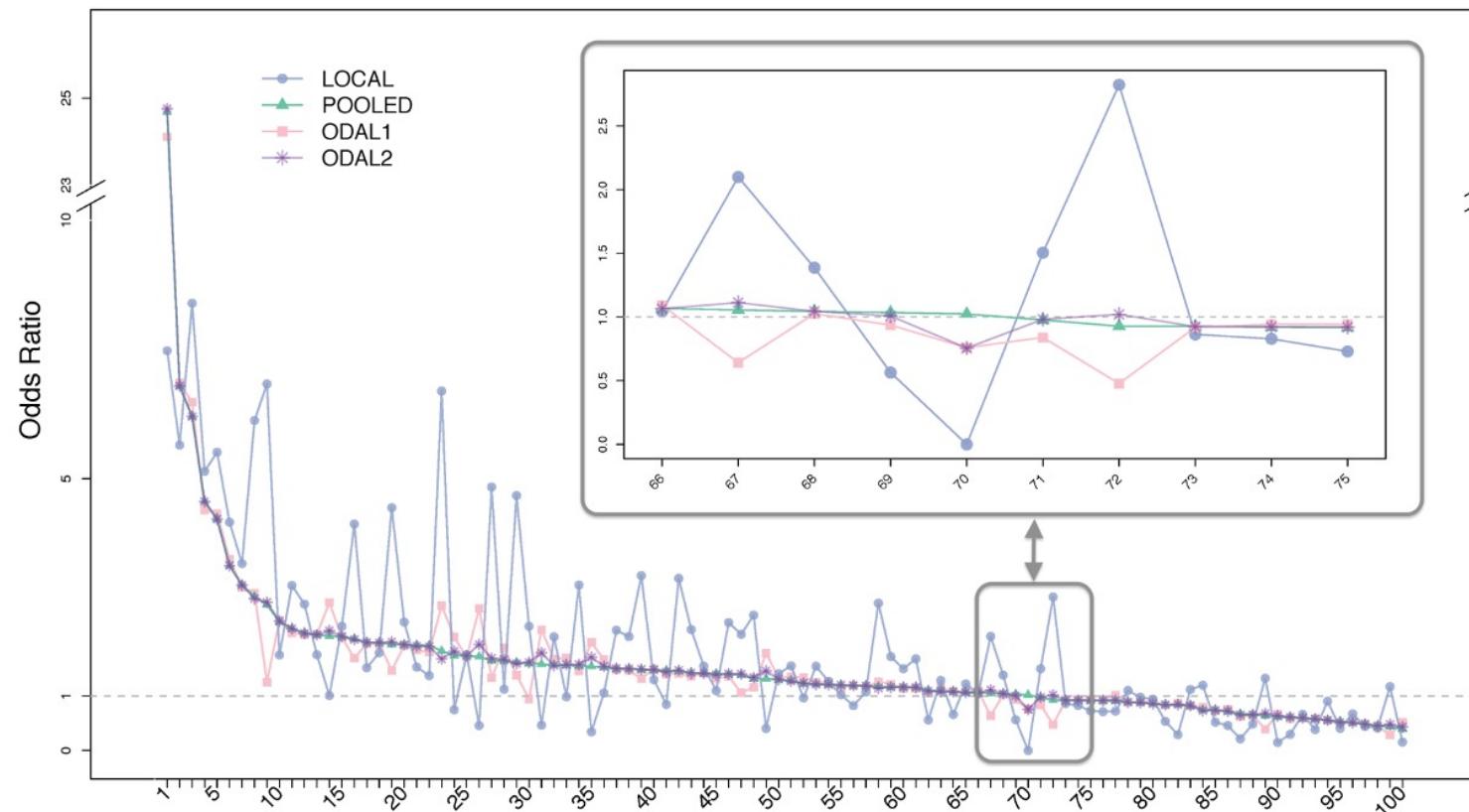
SHAKESPEARE LSTM, 54% ACCURACY

LSTM	E	B	u	IID	NON-IID
FEDSGD	1	∞	1.0	2488	3906
FEDAVG	1	50	1.5	1635 (1.5x)	549 (7.1x)
FEDAVG	5	∞	5.0	613 (4.1x)	597 (6.5x)
FEDAVG	1	10	7.4	460 (5.4x)	164 (23.8x)
FEDAVG	5	50	7.4	401 (6.2x)	152 (25.7x)
FEDAVG	5	10	37.1	192 (13.0x)	41 (95.3x)



Local-updating (FedAvg) can reduce communication rounds by ~100x relative to SGD

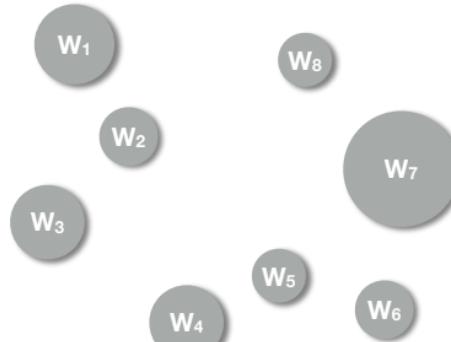
HIGHER ORDER DERIVATIVES



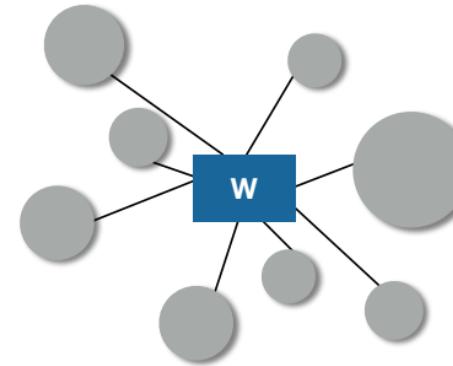
- ▶ ODAL2 needs to transfer an extra hessian matrix.
- ▶ ODAL2 provides more accurate estimation than ODAL1.

HOWTO MODEL HETEROGENOUS DATA

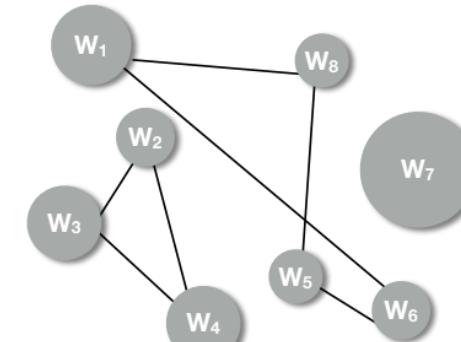
local



global



??



✓ personalized models
✗ don't learn from peers

✗ non-personalized models
✓ learn from peers

✓ personalized models
✓ learn from peers

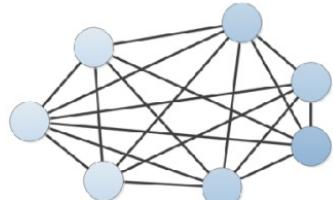
FEDERATED MULTI-TASK LEARNING

multi-task learning

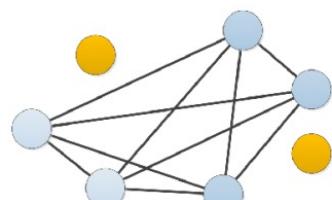
$$\min_{\mathbf{W}, \Omega} \sum_{t=1}^m \sum_{i=1}^{n_t} \ell_t(\mathbf{w}_t, \mathbf{x}_t^i) + \mathcal{R}(\mathbf{W}, \Omega)$$

models task relationship losses regularizer

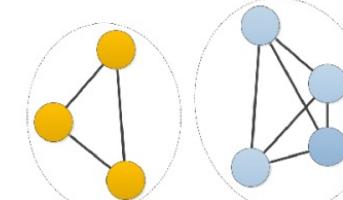
all tasks related



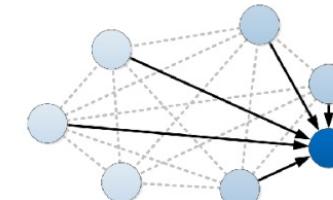
outlier tasks



clusters / groups

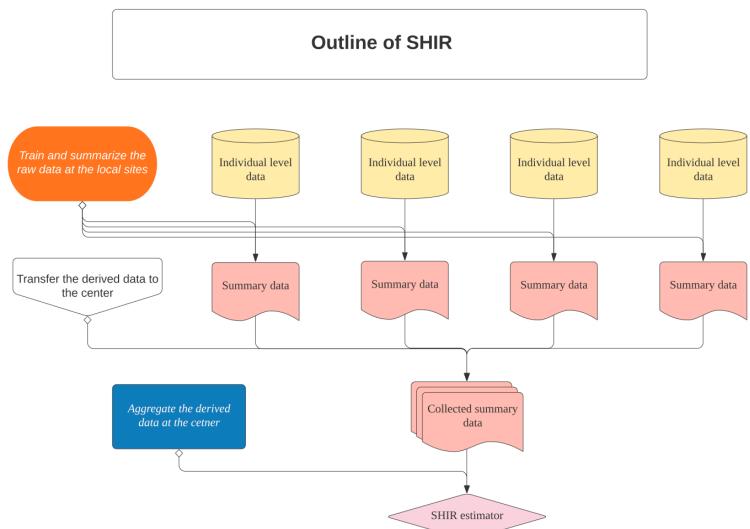


asymmetry



[EP, KDD 2004; JVB, NIPS 2009; ZCY, NIPS 2011]

AN EXAMPLE: SHIR



At the center node, we decompose the coefficients into mean effects and heterogeneous effects:

$$\{\mu, \alpha^{(\bullet)}\} \text{ where } \alpha^{(\bullet)} = (\alpha^{(1)}, \dots, \alpha^{(M)}), \beta^{(m)} = \mu + \alpha^{(m)}, \alpha^{(1)} + \dots + \alpha^{(M)} = 0,$$

and fit the quadratic-approximated integrative regression:

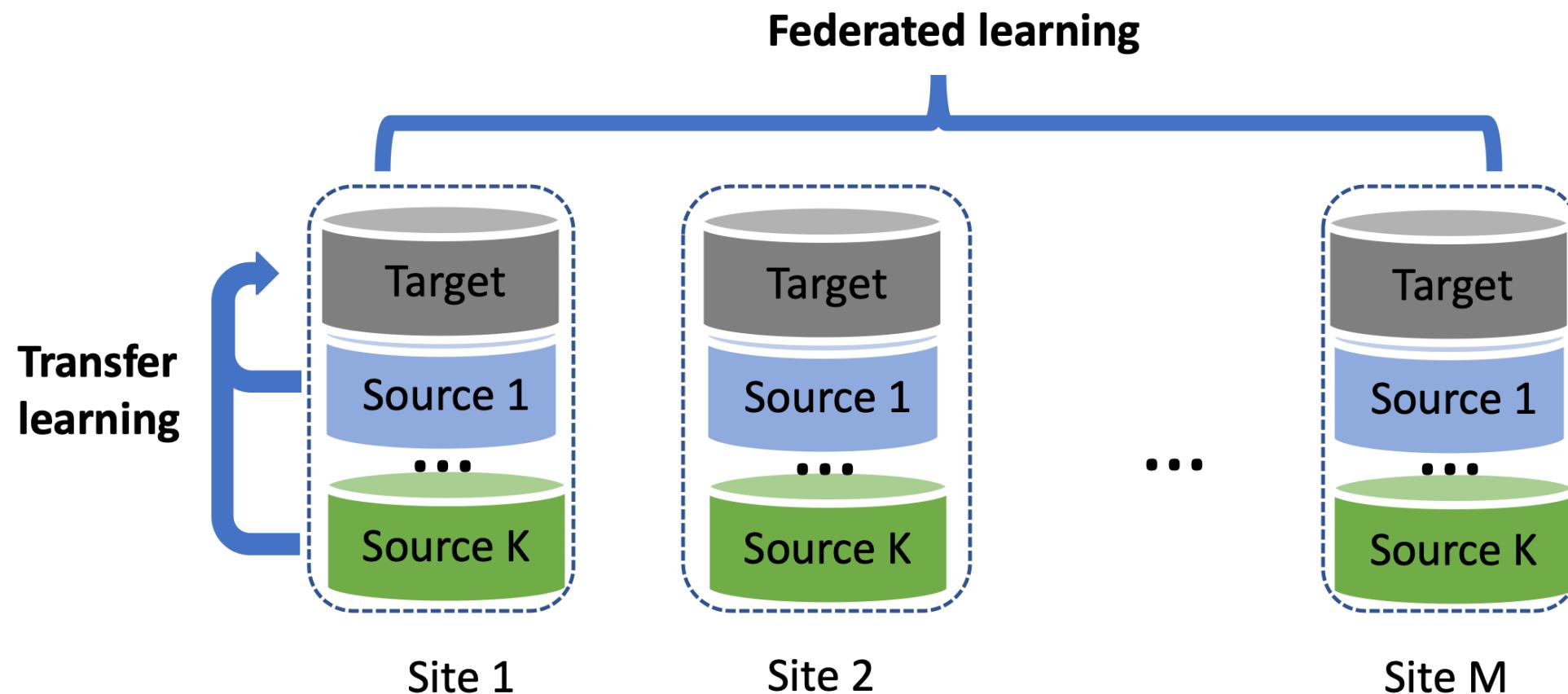
$$\{\hat{\mu}_{\text{SHIR}}, \hat{\alpha}_{\text{SHIR}}^{(\bullet)}\} = \operatorname{argmin}_{\mu, \alpha^{(\bullet)}} N^{-1} \sum_{m=1}^M n_m Q_m(\mu, \alpha^{(m)}) + \lambda (\|\mu - \mu_0\|_1 + \lambda_g \|\alpha_{-1}^{(\bullet)}\|_{2,1}),$$

$$\text{where } Q_m(\mu, \alpha^{(m)}) = (\mu + \alpha^{(m)})^\top \hat{H}_m(\mu + \alpha^{(m)}) - 2\hat{g}_m^\top(\mu + \alpha^{(m)}),$$

to obtain the SHIR estimator. Please see more details from the SHIR paper linked in the citation section.

A FEDERATED TRANSFER LEARNING APPROACH

- Li S, Cai T, Duan R. Targeting underrepresented populations in precision medicine: A federated transfer learning approach. *AOAS in press.*

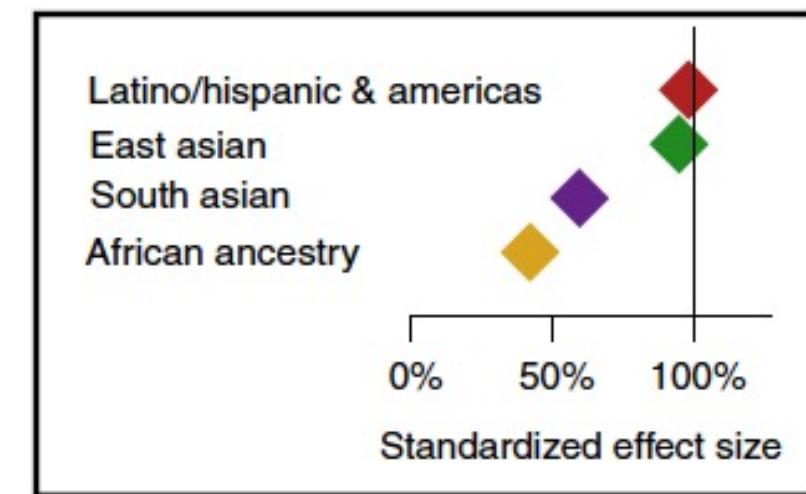


LACK OF REPRESENTATION OF BIOBANK DATA

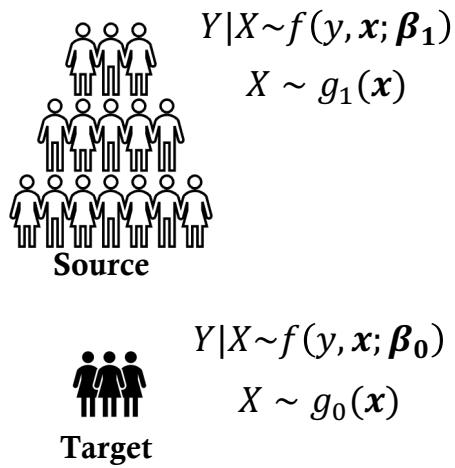
Duncan, Laramie, et al. *Nature communications* 10.1 (2019): 1-9.

Table 1: Sample size information of several large biobanks.

Study	N	%EA	%AA	%Hisp	%Asian
Million Veteran Program (MVP)†	825K	73	19	8	0
UK Biobank (UKB)†	500K	94	1.6	0.6	2.3
eMERGE†	83K	76	18	5	1
Mass General Brigham (MGB)†	112K	82	6	5	4
Penn Medicine Biobank (PMBB)†	62K	68	21	2.5	2.4
TOPMed	144K	40	32	16	10
All of us	400K	50	18	18	7



MODEL HETEROGENEITY



- ▶ Account for covariate shift (Pan & Yang 2010), i.e., $g_0(x) \neq g_1(x)$
 - Density ratio model: $w(x) = \frac{g_0(x)}{g_1(x)}$
 - E.g., exponential tilting model (Liang & Qin, 2000, Luo & Tsai 2012)

$$w(x; \eta) = \exp\{\eta_0 + \eta^T \psi(x)\}$$

- ▶ Outcome model shift , i.e., $\beta_0 \neq \beta_1$.
 - Specify common and population specific parameters (Duan et al, 2022)

$$\beta_0 = (\theta, \gamma_0); \beta_1 = (\theta, \gamma_1)$$

- Similarity constraints (Li et al, 2021, 2022)

$$d(\beta_0, \beta_1) < h$$

CENTRALIZED SETTING: WITHOUT DATA-SHARING CONCERNS

- ▶ Fitting a density ratio model:
 - Exponential titling is equivalent as a logistic regression model with outcome being the population indicator.
 - $Z = 1$ (target) $Z=0$ (source)
 - $P(Z=1 | X) = \text{expit}\{\eta_0 + \boldsymbol{\eta}^T \boldsymbol{\psi}(x)\}$
- ▶ Joint estimation of outcome model parameters

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} - \sum_{i \in I_T} \log f(x_i, y_i, \boldsymbol{\beta}_0) - \sum_{k=1}^K \sum_{i \in I_s} \widehat{w}_i \log f(x_i, y_i, \boldsymbol{\beta}_1) + \lambda \mathcal{P}(\boldsymbol{\beta}_0, \boldsymbol{\beta}_1)$$

- E.g., $\mathcal{P}(\boldsymbol{\beta}_0, \boldsymbol{\beta}_1) = \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_1\|_q$

FEDERATED SETTING

- ▶ Estimating density ratio model when source and target data cannot be pooled together

- For any $h(X)$, we have

$$E_T\{h(\mathbf{X})\} = E_s\{h(\mathbf{X})w(\mathbf{X}; \boldsymbol{\eta})\}$$

- We have the estimating equation for solving $\hat{\boldsymbol{\eta}}$

$$\frac{1}{n_T} \sum_{i \in I_T} h(\mathbf{x}_i) = \frac{1}{n_S} \sum_{i \in I_S} h(\mathbf{x}_i) w(\mathbf{x}_i; \boldsymbol{\eta})$$

- ▶ Approximation of the loss function: for any loss function $\mathbf{L}(\mathbf{b})$, we have the quadratic approximation with an initial value $\bar{\mathbf{b}}$

$$\tilde{\mathbf{L}}(\mathbf{b}; \bar{\mathbf{b}}) \propto \nabla \mathbf{L}(\bar{\mathbf{b}})^T (\mathbf{b} - \bar{\mathbf{b}}) + \frac{1}{2} (\mathbf{b} - \bar{\mathbf{b}})^T \nabla^2 \mathbf{L}(\bar{\mathbf{b}}) (\mathbf{b} - \bar{\mathbf{b}})$$

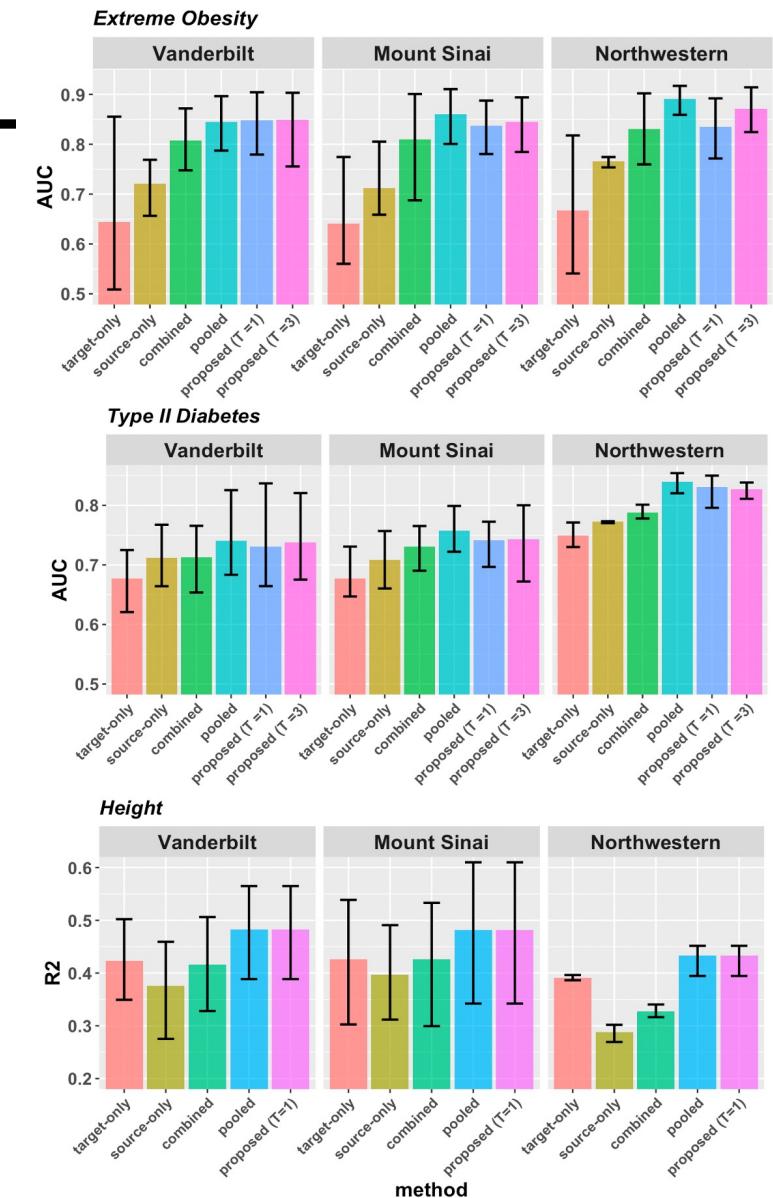
- ▶ $\frac{1}{n_T} \sum_{i \in I_T} h(\mathbf{x}_i), \nabla \mathbf{L}(\bar{\mathbf{b}}), \nabla^2 \mathbf{L}(\bar{\mathbf{b}})$ are summary statistics can be shared across sites
-

APPLICATION TO eMERGE

- Genetic risk models
- Integrate data from 7 sites at eMERGE
- Treat AA as target
- Three testing datasets from Vanderbilt, Mt Sinai, and Northwestern

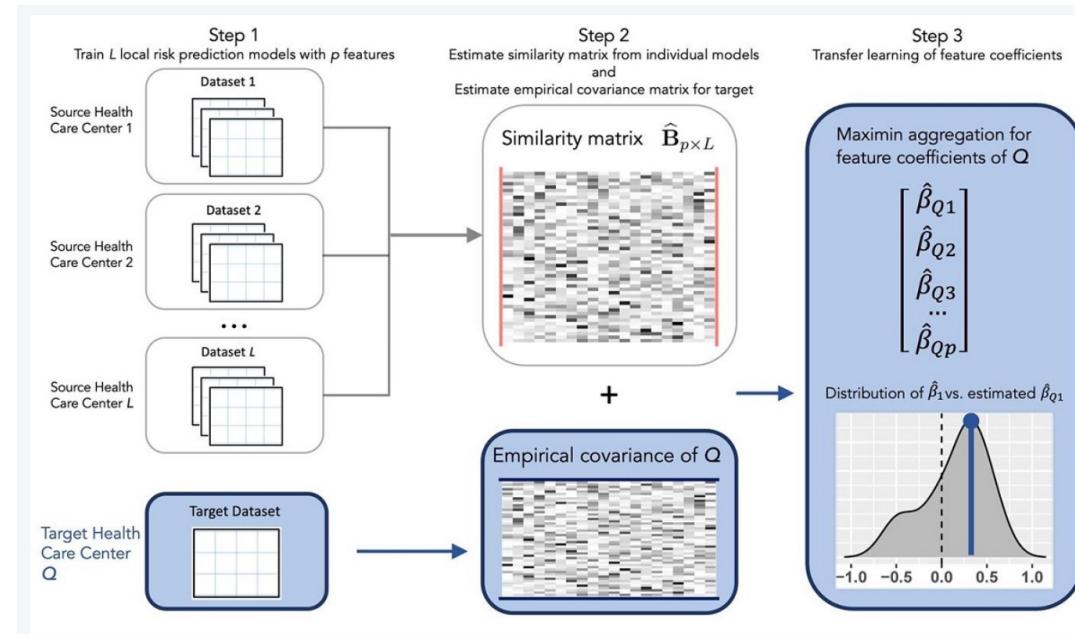
Table 1: Sample sizes of the target (AA) and the source (all other races) populations for type II diabetes, extreme obesity and height across seven sites.

Site Name	Extreme Obesity		Height		Type II Diabetes	
	source	target	source	target	source	target
Geisinger Health System	1023	4	2643	8	3081	9
Group Health Cooperative	273	13	2710	94	262	16
Marshfield Clinic	1497	0	3959	2	3977	3
Mayo Clinic	1320	3	4790	11	2880	10
Mount Sinai	534	1142	524	3402	606	3653
Northwestern University	1305	167	1178	180	1223	301
Vanderbilt University	1523	646	6598	1852	2580	1646



ZERO-SHOT LEARNING: INCORPORATING PRE-TRAINED MODELS

- Model aggregation
- Transfer learning
- Maximin Learning

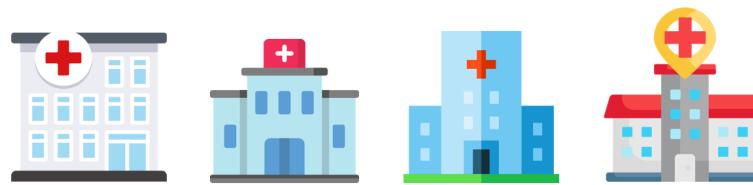


<https://shiny.parse-health.org/ARRTLE/>

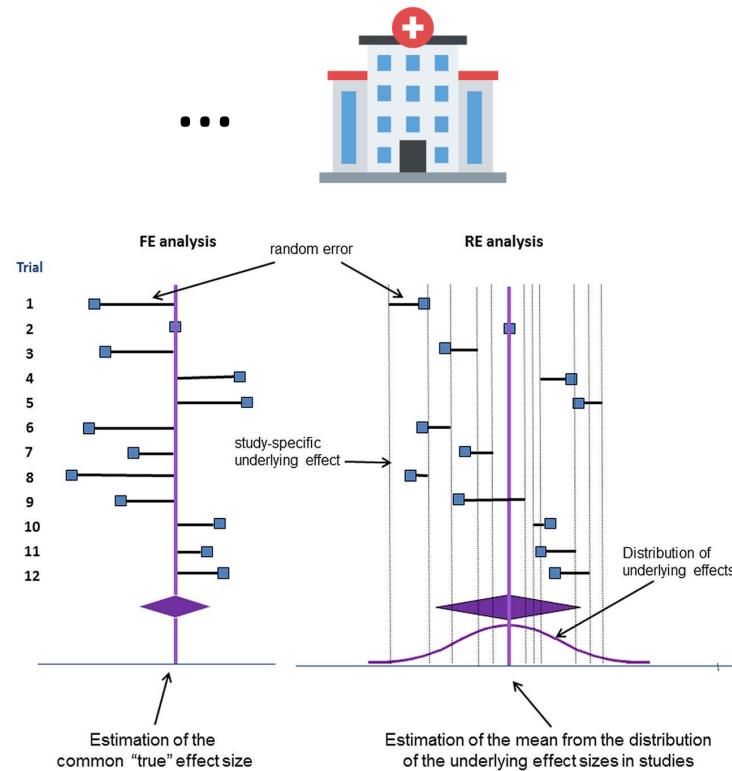
- Lecué, Guillaume, and Philippe Rigollet. "Optimal learning with Q-aggregation." (2014): 211-224.
- Parisi, Fabio, et al. "Ranking and combining multiple predictors without labeled data." *Proceedings of the National Academy of Sciences* 111.4 (2014): 1253-1258.
- Gu, Tian, Yi Han, and Rui Duan. "Robust angle-based transfer learning in high dimensions." *arXiv preprint arXiv:2210.12759* (2022).
- Wang, Zhenyu, Peter Bühlmann, and Zijian Guo. "Distributionally robust machine learning with multi-source data." *arXiv preprint arXiv:2309.02211* (2023).

SECTION III: FEDERATED STATISTICAL INFERENCE

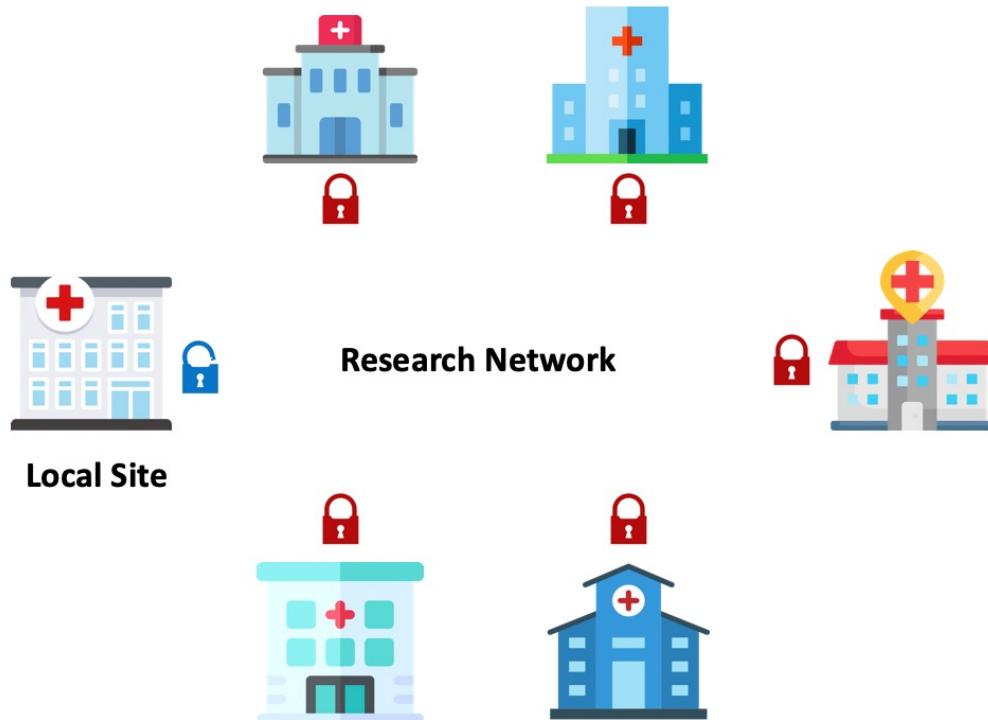
STATISTICAL INFERENCE OF SHARED PARAMETERS ACROSS SITES



- Benchmark: Meta-analysis
- Average across local estimates



IID DATA: COMMUNICATION-EFFICIENT SURROGATE LIKELIHOOD



- Patient-level data from one site is accessible
- Only aggregated information can be shared between site

- ▶ **Combined likelihood function** (if data could be shared)

$$L(\beta) = \frac{1}{nK} \sum_{j=1}^K \log f(d_{ij}; \beta)$$

- ▶ **Local likelihood function** (assume local site to be the first site, j=1)

$$L_1(\beta) = \frac{1}{n} \sum_{i=1}^n \log f(d_{i1}; \beta)$$

- ▶ For an initial value $\bar{\beta}$,

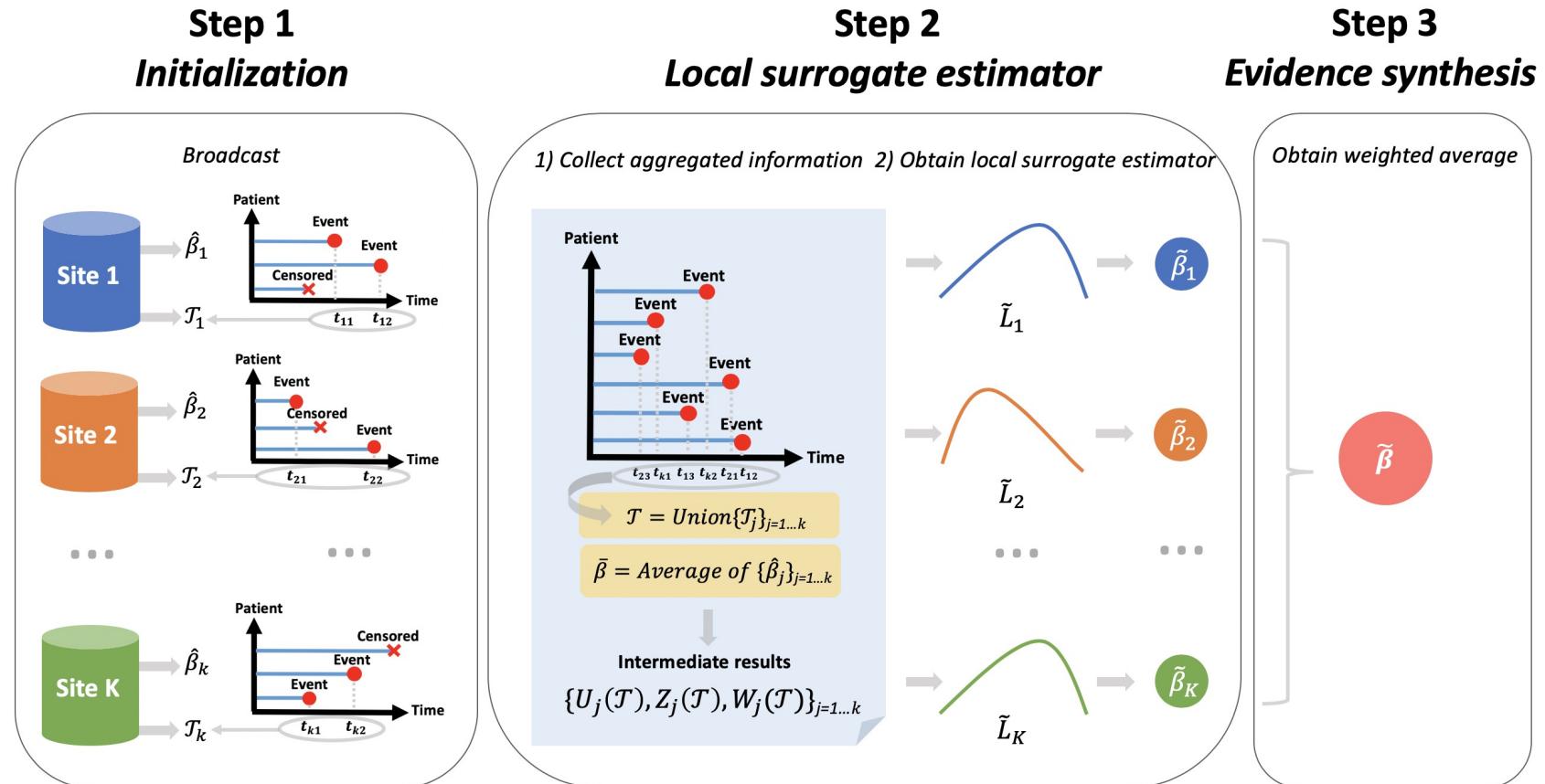
$$L(\beta) = L(\bar{\beta}) + \nabla L(\bar{\beta})^T (\beta - \bar{\beta}) + \sum_{t=1}^{\infty} \frac{1}{t!} \nabla^t L(\bar{\beta}) (\beta - \bar{\beta})^{\otimes t}$$

$$\begin{aligned} L_1(\beta) &= L_1(\bar{\beta}) + \nabla L_1(\bar{\beta})^T (\beta - \bar{\beta}) + \sum_{t=1}^{\infty} \frac{1}{t!} \nabla^t L_1(\bar{\beta}) (\beta - \bar{\beta})^{\otimes t} \\ \sum_{t=1}^{\infty} \frac{1}{t!} \nabla^t L_1(\bar{\beta}) (\beta - \bar{\beta})^{\otimes t} &= L_1(\beta) - L_1(\bar{\beta}) - \nabla L_1(\bar{\beta})^T (\beta - \bar{\beta}) \end{aligned}$$

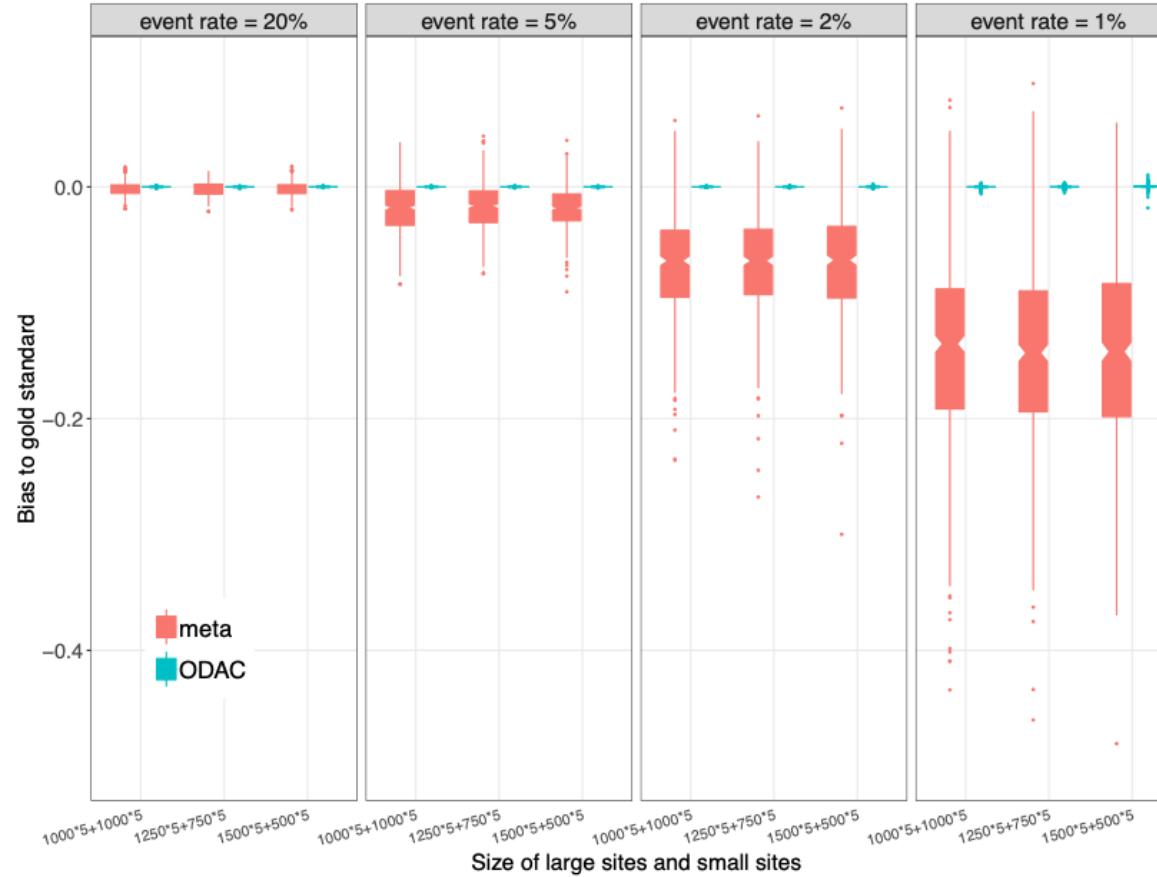
SL function

$$\tilde{L}^1(\beta) = L_1(\beta) + \{\nabla L(\bar{\beta}) - \nabla L_1(\bar{\beta})\}^T \beta$$

ODAC



BENEFIT IN STUDYING RARE EVENTS



- ▶ Meta-analysis has increasing bias when event is rarer.
- ▶ ODAC provides estimates close to the pooled analysis.

ALLOW SITE-SPECIFIC PARAMETERS

$$Y_{ij} \sim f(y; \boldsymbol{\beta}, \boldsymbol{\gamma}_j)$$

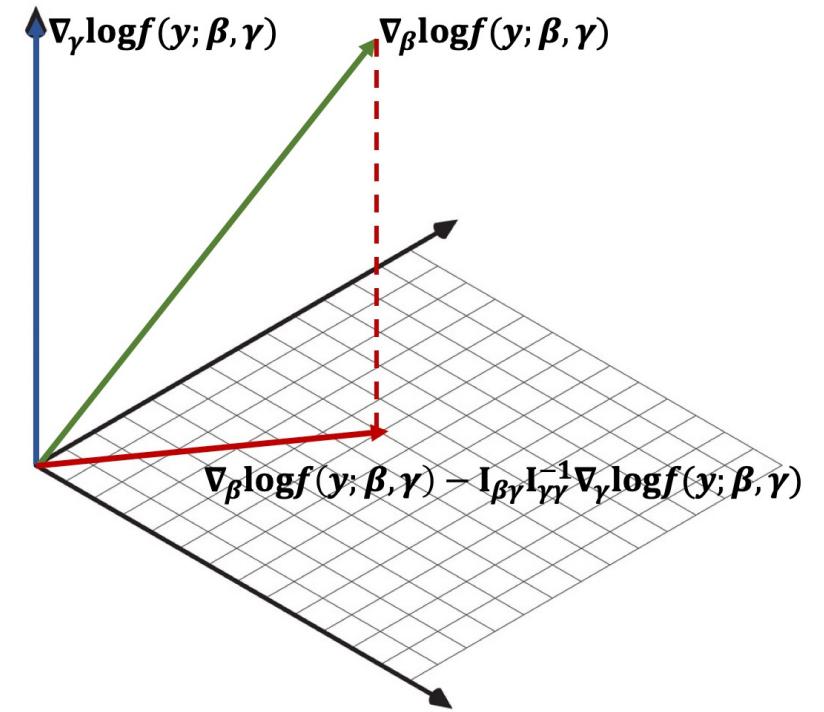
- β is the *parameter of interest*
 - γ_j is the *site-specific nuisance parameter*--- allow site to be a covariate variable, allow interaction terms between site and other covariates.
-

Change the target function to the combined efficient score function

$$L(\beta; \bar{\Gamma}) = \frac{1}{Kn} \sum_{j=1}^K \sum_{i=1}^n \log f(y_{ij}; \beta, \bar{\gamma}_j)$$



$$\{ S(\beta; \bar{\Gamma}) = \frac{1}{Kn} \sum_{j=1}^K \sum_{i=1}^n \left\{ \nabla_\beta \log f(y_{ij}; \beta, \bar{\gamma}_j) - I_{\beta\gamma}^{(j)} I_{\gamma\gamma}^{(j)-1} \nabla_\gamma \log f(y_{ij}; \beta, \bar{\gamma}_j) \right\}$$



Find a local function and construct higher-order match

- We want to find some function $g(y; \beta)$, such that

$$\mathbb{E}_{f_1} \{\nabla^t g(Y_{i1}; \beta)\} = \mathbb{E} \{\nabla_\beta^t S(\beta; \bar{\Gamma})\} = \frac{1}{K} \sum_{j=1}^K \mathbb{E}_{f_j} \nabla_\beta^t \left\{ \nabla_\beta \log f(Y_{ij}; \beta, \bar{\gamma}_j) - I_{\beta\gamma}^{(j)} I_{\gamma\gamma}^{(j)-1} \nabla_\gamma \log f(Y_{ij}; \beta, \bar{\gamma}_j) \right\}$$

- Density ratio tilting
- For initial estimator $\bar{\beta}, \bar{\gamma}_j$, define

$$g(y; \beta) = \frac{1}{K} \sum_{j=1}^K \frac{\mathbf{f}(y; \bar{\beta}, \bar{\gamma}_j)}{\mathbf{f}(y; \bar{\beta}, \bar{\gamma}_j)} \left\{ \nabla_\beta \log f(y; \beta, \bar{\gamma}_j) - \tilde{H}_{\beta\gamma}^{(j)} \tilde{H}_{\gamma\gamma}^{(j)-1} \nabla_\gamma \log f(y; \beta, \bar{\gamma}_j) \right\}$$

- Define

Surrogate efficient score function

$$\tilde{U}(\beta) = \mathbf{s}(\bar{\beta}; \bar{F}) + \mathbf{U}_1(\beta) - \mathbf{U}_1(\bar{\beta})$$

where $U_1(\beta) = \sum_{i=1}^n g(y_{i1}; \beta)/n$

- Obtain the proposed estimator by $\tilde{U}(\beta) = 0$

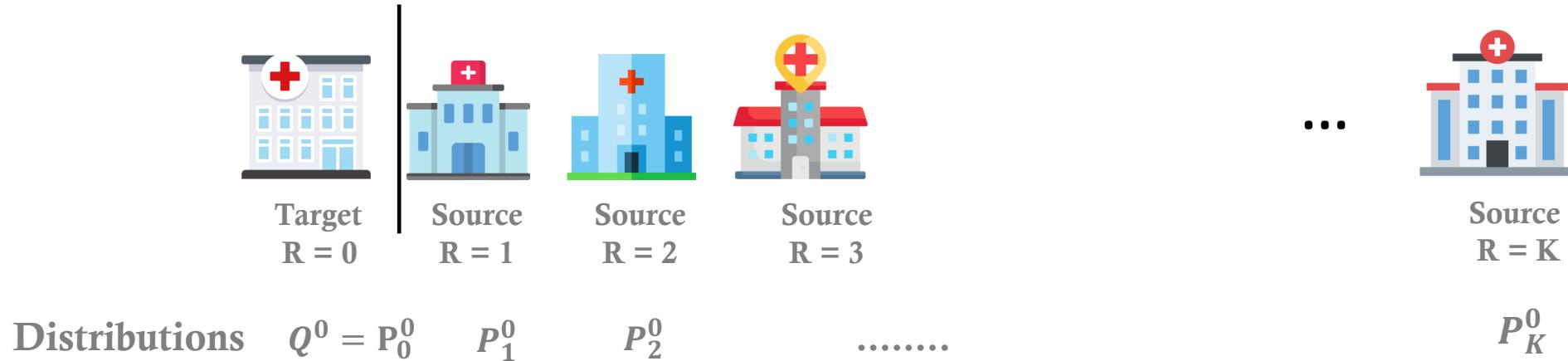
- Each site need to provide $\bar{\gamma}_j$, and

$$S_j(\bar{\beta}; \bar{\gamma}_j) = \frac{1}{n} \sum_{i=1}^n \left\{ \nabla_{\beta} \log f(Y_{ij}; \bar{\beta}, \bar{\gamma}_j) - H_{\beta\gamma}^{(j)} H_{\gamma\gamma}^{(j)-1} \nabla_{\gamma} \log f(Y_{ij}; \bar{\beta}, \bar{\gamma}_j) \right\}$$

THEORETICAL PROPERTIES

Desired Properties	Meta-analysis	Proposed estimator (T=1)	Proposed estimator (T=2)
Consistency	consistent	consistent	consistent
Distance to the Gold Standard Estimator	$\frac{C}{\sqrt{Kn}}$	$\leq \frac{C}{n}$	$\leq \frac{C}{n\sqrt{K}} + \frac{C}{n\sqrt{n}}$
Asymptotic Normality	asymptotic normal	asymptotic normal	asymptotic normal
Asymptotic Efficient	not efficient	efficient	efficient

STATISTICAL INFERENCE IN A TARGET POPULATION



Estimand of interest is defined in the target $\psi(Q^0)$

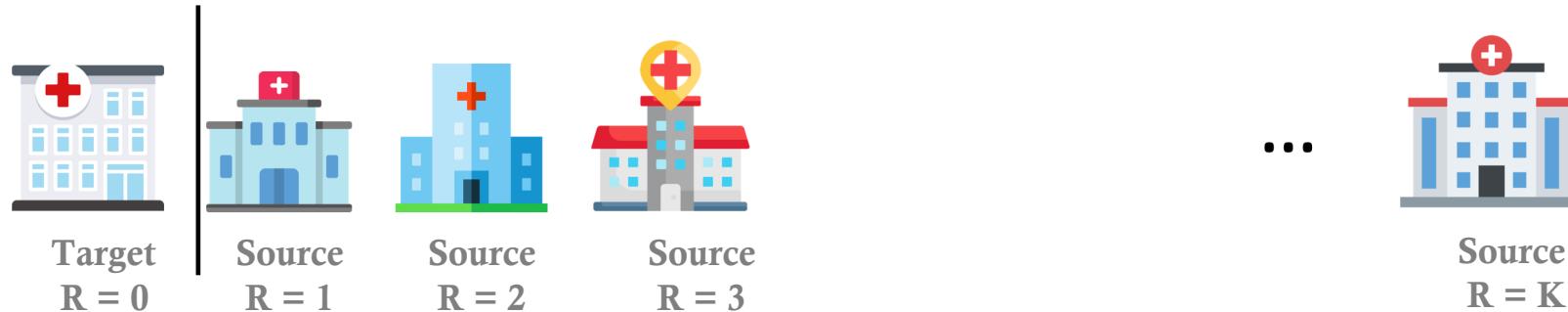
How can all sites work collaboratively to estimate $\psi(Q^0)$?

CAUSAL INFERENCE OF TREATMENT EFFECTS



Consortium for Clinical Characterization of COVID-19 by EHR

ESTIMATE TARGET AVERAGE TREATMENT EFFECT (TATE)



- ▶ $Y^{(a)}$: potential outcome of patients under treatment $a = 0, 1$
- ▶ Estimand: $\Delta_T = \mathbf{E}(Y^{(1)} - Y^{(0)} | R = 0) = \mathbf{E}_{Q^0}(Y^{(1)} - Y^{(0)})$
- ▶ At target, we assume positivity, consistency, no unmeasured confounding assumptions for the identification of TATE
$$\mathbf{E}_{Q^0;A,X}(\mathbf{E}_{Q^0;Y|A,X}(Y|A = 1, X)) - \mathbf{E}_{Q^0;A,X}(\mathbf{E}_{Q^0;Y|A,X}(Y|A = 0, X))$$
- ▶ Distribution shifts in $[Y, A, X]$ across sites

HOW CAN A HETEROGENEOUS SOURCE SITE HELP WITH ESTIMATING THE TARGET ESTIMAND?

- ▶ Distributional shifts in $[A, X]$ can be handled by adjusting sample weights according to the density ratio and propensity.
- ▶ The conditional distribution $[Y | A, X]$ is the key for consistently estimating the target ATE.
- ▶ Many assumes exchangeability: $[Y | A, X]$ is the same across sites
 - Causal transportability: Bareinboim & Pearl, 2014; Rudolph & van der Laan, 2017; Dahabreh et al., 2019; Rudolph & van der Laan, 2017; Dahabreh & Hernán, 2019
 - More general data fusion: Athey et al., 2019; Kallus et al., 2020; Lee et al., 2023; Brantner et al., 2023

WHAT IS MORE PRACTICAL IN THE REAL-WORLD SETTINGS?



- Setting 1:
 - Among K source sites, exchangeability holds only among a subset (could be empty) of all source sites.
 - We do not know the membership of the subset

Federated Adaptive Causal Estimation (FACE) of Target Treatment Effects

Larry Han^{1,2}, Jue Hou³, Kelly Cho⁴, Rui Duan^{1†}, Tianxi Cai^{1,5†}

1 Department of Biostatistics, Harvard University

2 Department of Health Sciences, Northeastern University

3 Division of Biostatistics, University of Minnesota

4 Massachusetts Veterans Epidemiology Research and
Information Center, US Department of Veteran Affairs

5 Department of Biomedical Informatics, Harvard Medical School

† Co-corresponding authors

Other recent work: Yang et al 2023—test and pool, etc

OUTLINE OF THE METHOD

- Target site: obtain $\widehat{\Delta}_{T,0}$ from augmented inverse probability weighting (AIPW) approach.

$$\begin{aligned}\widehat{\Delta}_{T,0} &= M_0 + \delta_0 \\ M_0 &= n^{-1} \sum_{i \in I_t} \{\widehat{m}(1, X_i) - \widehat{m}(0, X_i)\} \\ \delta_0 &= n^{-1} \sum_{i \in I_T} \left\{ \frac{I(A_i=1)}{\widehat{\pi}_0(1, X_i)} - \frac{I(A_i=0)}{\widehat{\pi}_0(0, X_i)} \right\} \{Y_i - \widehat{m}(A_i, X_i)\}\end{aligned}$$

- Source site(s): obtain $\widehat{\Delta}_{T,k}$ with the help of a density ratio model

$$w_k(\mathbf{X}) = \frac{f(\mathbf{X}|R=0)}{f(\mathbf{X}|R=k)}$$

$$\delta_k = n_k^{-1} \sum_{i \in I_k} \widehat{w}_k(X_i) \left\{ \frac{I(A_i=1)}{\widehat{\pi}_k(1, X_i)} - \frac{I(A_i=0)}{\widehat{\pi}_k(0, X_i)} \right\} \{Y_i - \widehat{m}(A_i, X_i)\}$$

$$\widehat{\Delta}_{T,k} = M_0 + \delta_k$$

- ▶ Combine source and target estimators through an adaptive aggregation step
-

ADAPTIVE AGGREGATION

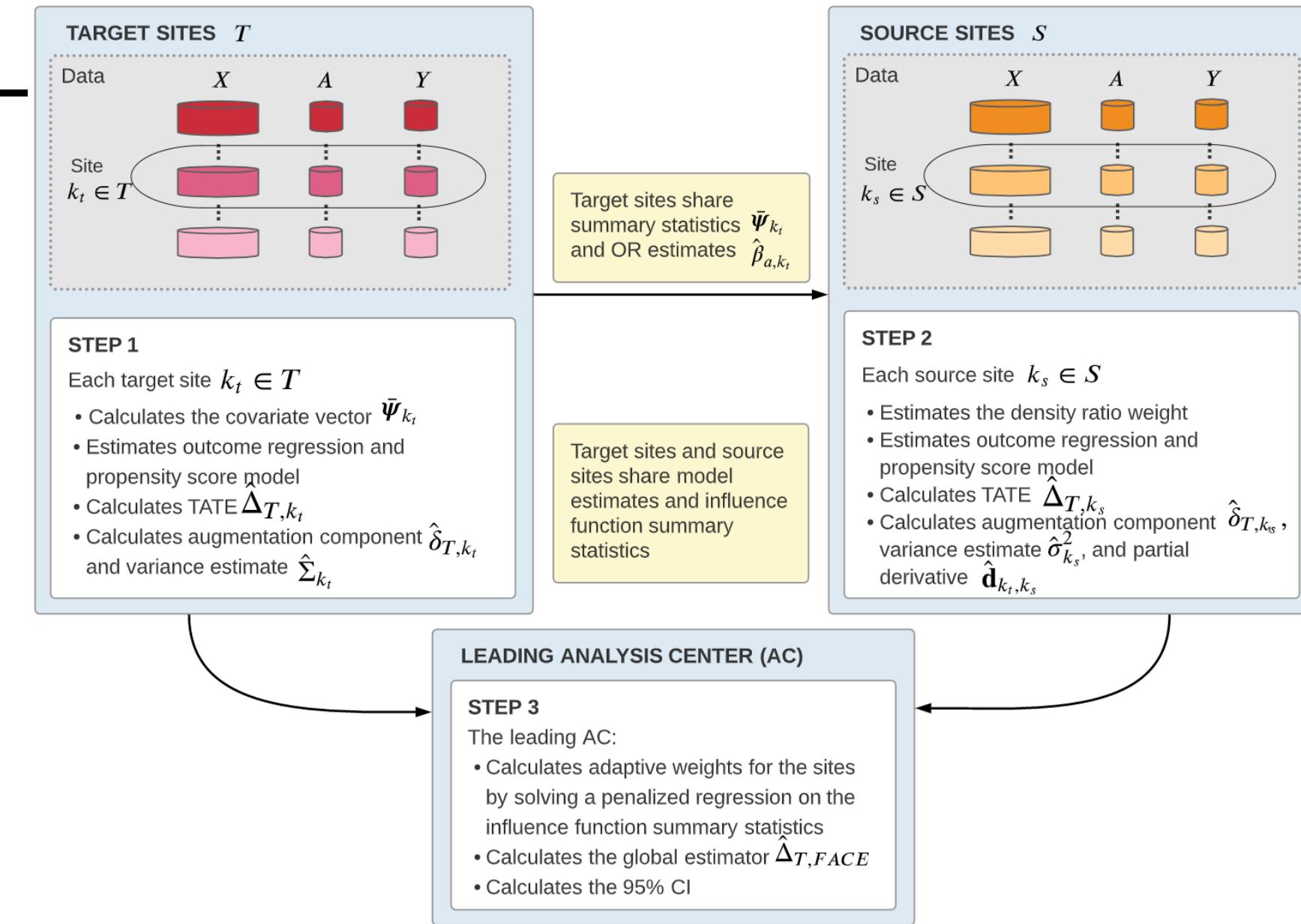
- ▶ Anchor and aggregation

$$\widehat{\Delta}_T = \widehat{\Delta}_{T,0} + \sum_{k \in S} \eta_k \{ \widehat{\Delta}_{T,k} - \widehat{\Delta}_{T,0} \}$$

- ▶ Best weights learned through

$$\widehat{\boldsymbol{\eta}} = \operatorname{argmin} \operatorname{Var}\{ \widehat{\Delta}_{T,0} + \sum_{k \in S} \eta_k \{ \widehat{\Delta}_{T,k} - \widehat{\Delta}_{T,0} \} \} + \lambda \sum_{k \in S} |\eta_k| (\widehat{\Delta}_{T,k} - \widehat{\Delta}_{T,0})^2$$

- ▶ Data splitting needed for adaptive aggregation
 - ▶ Theoretical guarantees for asymptotical normality when target OR or PS model is correctly specified.
 - ▶ Efficiency gain when exchangeability holds in at least one site
-



APPLICATION TO COMPARATIVE EFFECTIVENESS OF COVID VACCINES

- To illustrate FACE, we study the comparative effectiveness of BNT162b2 (Pfizer) versus mRNA-1273 (Moderna) for the prevention of COVID-19 outcomes in five VA sites.
- It is of interest to understand the real-world effectiveness of these vaccines, but head-to-head comparisons have been rare.
- Inclusion criteria:
 - ▶ veteran status,
 - ▶ ≥ 18 years by Jan 1, 2021,
 - ▶ no previously documented COVID-19 infection, and
 - ▶ documented two-dose COVID-19 vaccination with either Pfizer or Moderna between Jan 1-Mar 24, 2021.
- All models were adjusted for age, sex, ethnicity, residence, comorbidities (CLD, CVD, HTN, T2DM, CKD, Autoimmune diseases, Obesity), and time.
- The outcomes of interest were documented SARS-CoV-2 infection either 120 or 180 days after baseline, and death with COVID-19 infection either 120 or 180 days after baseline.

APPLICATION TO COMPARATIVE EFFECTIVENESS OF VACCINES

- Han, Larry, et al. (2023+) "Federated Adaptive Causal Estimation (FACE) of Target Treatment Effects." JASA minor revision. *arXiv:2112.09313*.

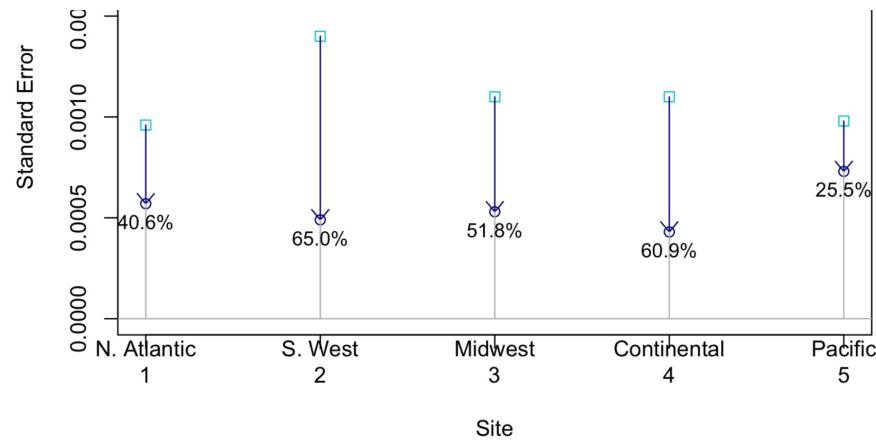
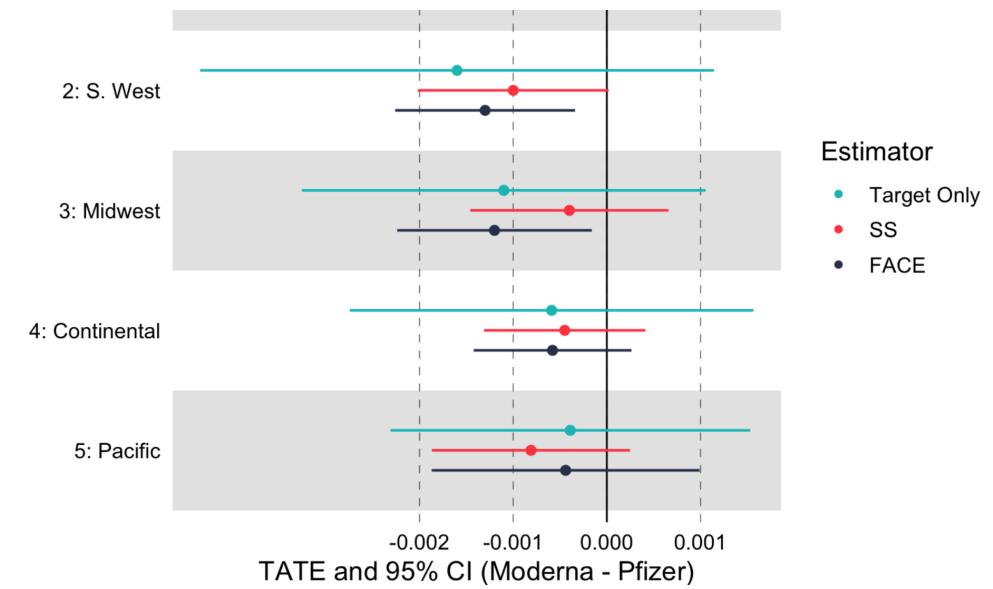


Figure: Reductions in SE are substantial using FACE



WHAT IS MORE PRACTICAL IN THE REAL-WORLD SETTINGS?

- When exchangeability does not hold in any of the source sites, the previous methods may be worse than the target only estimator.
- It does not distinguish two sites, one with a small shift in $Y | A, X$, and one with a large shift

NEW PROPOSAL

Our more recent proposal:

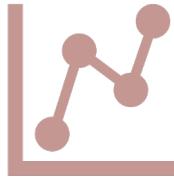
- ▶ A semiparametric model for the shifts in $Y|A, X$ (selection bias model)

$$p_k^0(y|x, a) = w_k^*(y, x, a; \beta_k, Q^0)q^0(y|x, a),$$

$$w_k^*(y, x, a; \beta_k, Q^0) = \frac{w_k(y, x, a; \beta_k)}{E_{Q^0}(w_k(y, x, a; \beta_k)|X = x, A = a)}$$

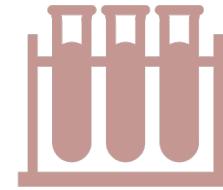
- ▶ For example, the exponential tilting model $w_k(y, x, a; \beta_k) = \exp(\phi(z)^T \beta)$
 - ▶ Efficient data integration through EIF
-

PROTECT AGAINST MODEL MISSPECIFICATION



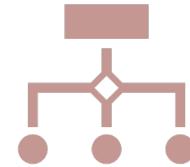
Use Domain Knowledge

Understand the meaning of the distributional shift parameters, e.g., shifts in log odds or risks in stratified groups.



Perform Goodness-of-Fit Tests and Sensitivity Analyses

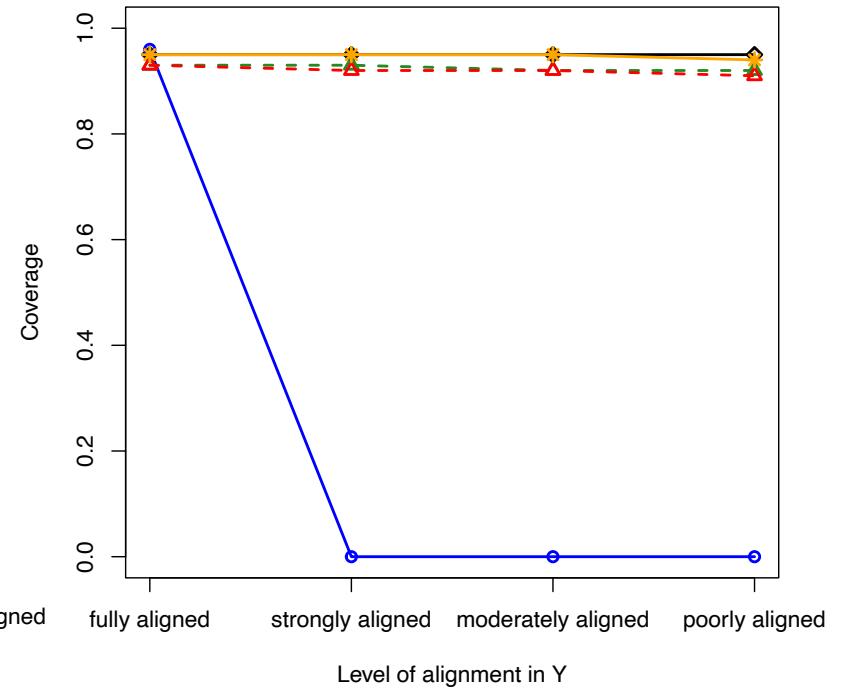
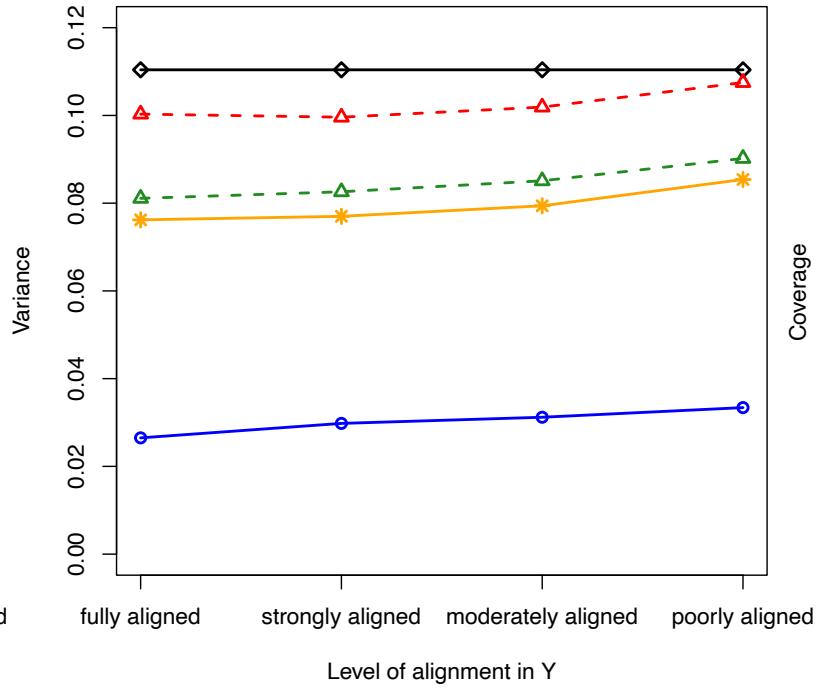
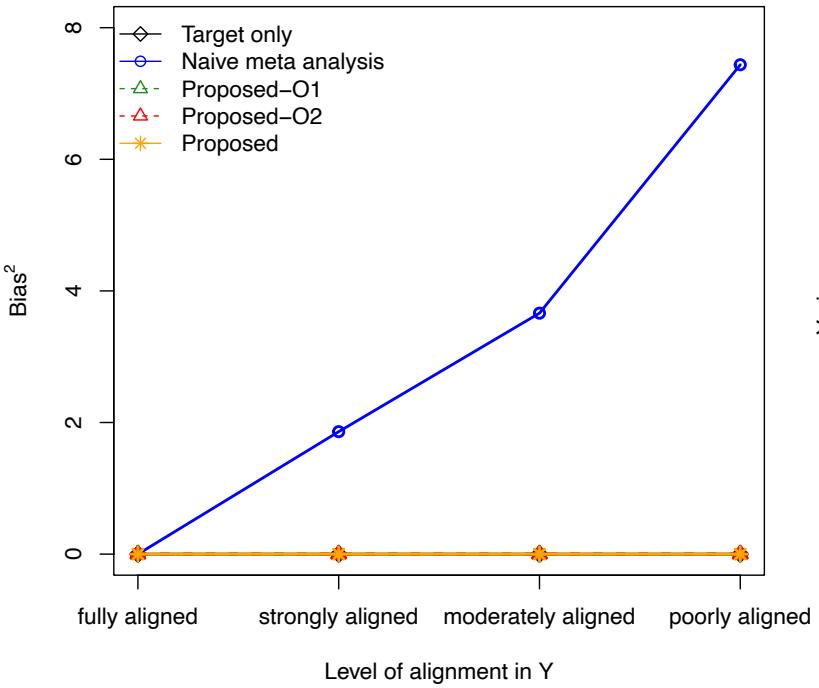
Regularly test the fit of the model to ensure accuracy. Conduct sensitivity analyses to assess the robustness of the model.



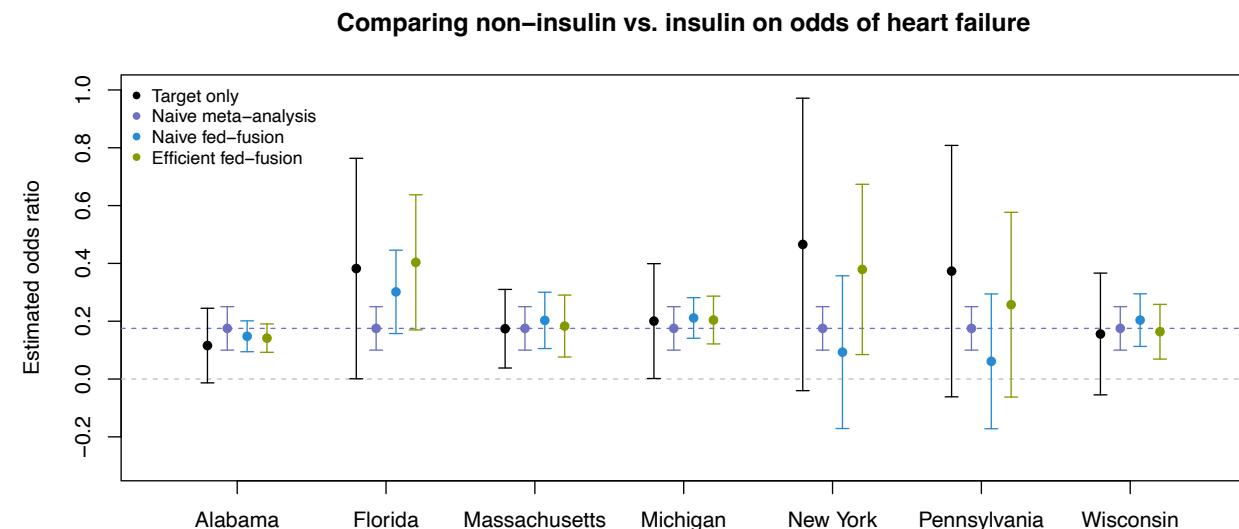
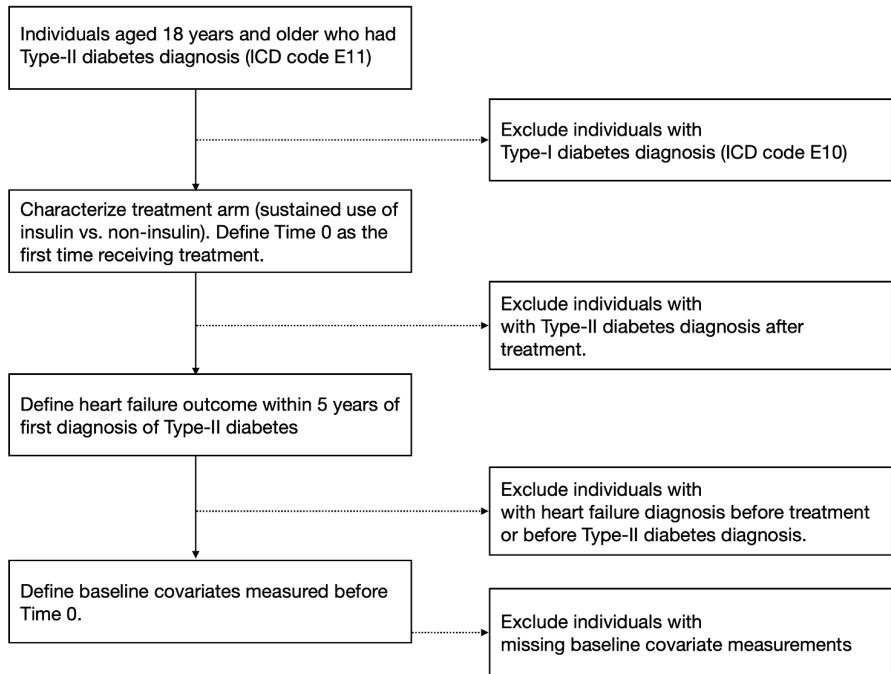
Over-parametrize the Model

Include more basis functions in the density ratio model---loss of efficiency

SIMULATION



REAL DATA—AOU RESEARCH HUB



Preprint coming shortly

DISCUSSION

- can we *effectively* and *efficiently* learn from heterogeneous sources?
 - communication-efficient distributed optimization, heterogeneity-aware optimization, personalization and multi-task learning, connections to pretraining
- can we provide *trustworthy* learning schemes?
 - techniques efficient/accurate/usable private and secure learning in large-scale distributed settings, threat modeling and auditing, robustness to attacks
- can we develop systems/tools that *support* collaboration?
 - game-theoretic perspectives on collaborative learning, data markets, data valuation and risk assessment, incorporating social dynamics
- Fairness: can we equalize performance across diverse networks?

REFERENCE

- Lecture Notes, Federated and Collaborative Learning, Fall 2023, Virginia Smith. <https://www.cs.cmu.edu/~smithv/10719/> Access date: 07-05-2024.
 - Li, Li, et al. "A review of applications in federated learning." *Computers & Industrial Engineering* 149 (2020): 106854.
 - Mammen, Priyanka Mary. "Federated learning: Opportunities and challenges." *arXiv preprint arXiv:2101.05428* (2021).
 - Zhang, Chen, et al. "A survey on federated learning." *Knowledge-Based Systems* 216 (2021): 106775.
 - Chen, Shuxiao, Bo Zhang, and Ting Ye. "Minimax rates and adaptivity in combining experimental and observational data." *arXiv preprint arXiv:2109.10522* (2021).
 - BAREINBOIM, E. & PEARL, J. (2014). Transportability from multiple environments with limited experiments: Completeness results. *Advances in neural information processing systems* 27, 280–288.
 - BRANTNER, C. L., CHANG, T.-H., NGUYEN, T. Q., HONG, H., DI STEFANO, L. & STUART, E. A. (2023). Methods for integrating trials and non-experimental data to examine treatment effect heterogeneity. *arXiv preprint arXiv:2302.13428*
 - DAHABREH, I. J. & HERNÁN, M. A. (2019). Extending inferences from a randomized trial to a target population. *Eur. J. Epidemiol.* 34, 719–722
 - LI, S. & LUEDTKE, A. (2023). Efficient estimation under data fusion. *Biometrika* 110, 1041–1054
 - RUDOLPH, K. E. & VAN DER LAAN, M. J. (2017). Robust estimation of encouragement-design intervention effects transported across sites. *J. R. Stat. Soc.* 79, 1509
 - HAN, L., HOU, J., CHO, K., DUAN, R. & CAI, T. (2021). Federated adaptive causal estimation (face) of target treatment effects. *arXiv preprint arXiv:2112.09313* .
 - YANG, S., GAO, C., ZENG, D. & WANG, X. (2023). Elastic integrative analysis of randomised trial and real-world data for treatment heterogeneity estimation. *Journal of the Royal Statistical Society Series B:Statistical Methodology* 85, 575–596.
-