

---

# Differentially private secure distributed logistic regression

---

**Xiaoqian Jiang**

Biomedical Informatics

University of California San Diego

**IPAM workshop: algorithmic challenges of protecting biomedical data**

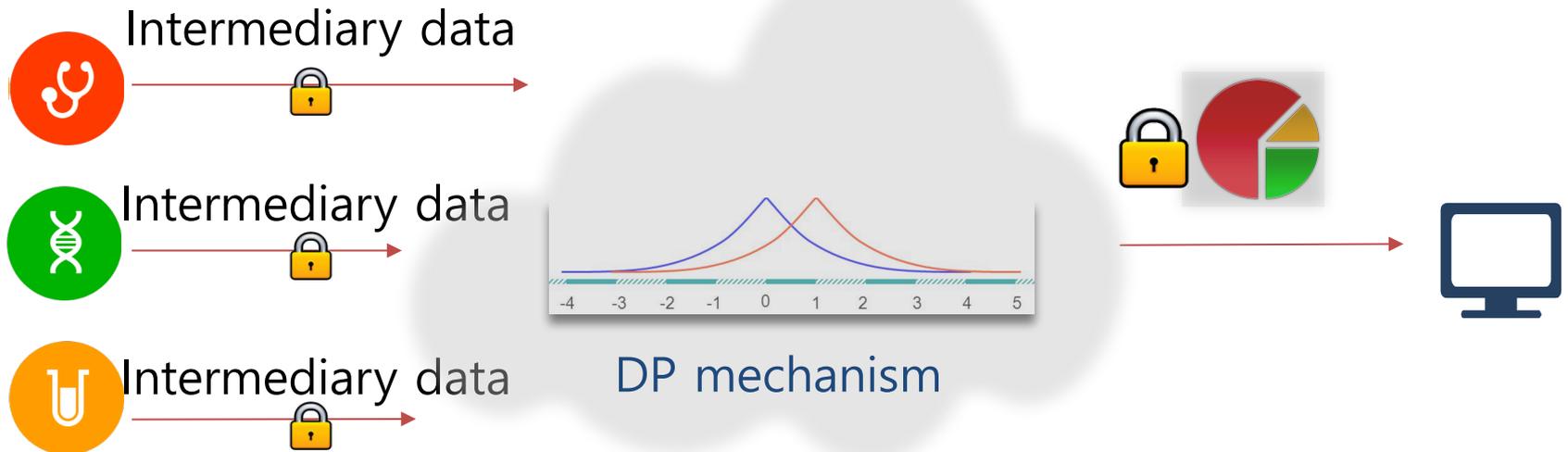
**1/11/2018**

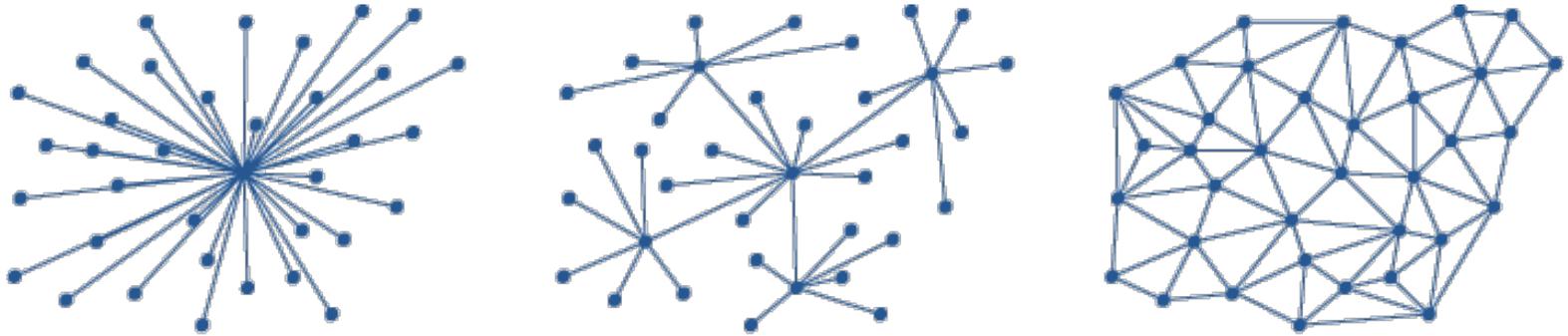


# Many attack models have been discovered...

- **Malin 2005:** Trails of hospital visit pattern might lead to information disclosure
- **Machanavajjhala 2007:** Demographic statistics for certain cohort can lead to privacy leakage.
- **Loukides 2010:** Distribution of disease can lead to re-identification
- **Sweeney 2014:** Demographics combined with phenotypes provide strong clues to reveal individuals' information
- **Bonomi 2017:** Hospital visit frequency and interval can lead to re-identification

# Homomorphic encryption and differential privacy might help

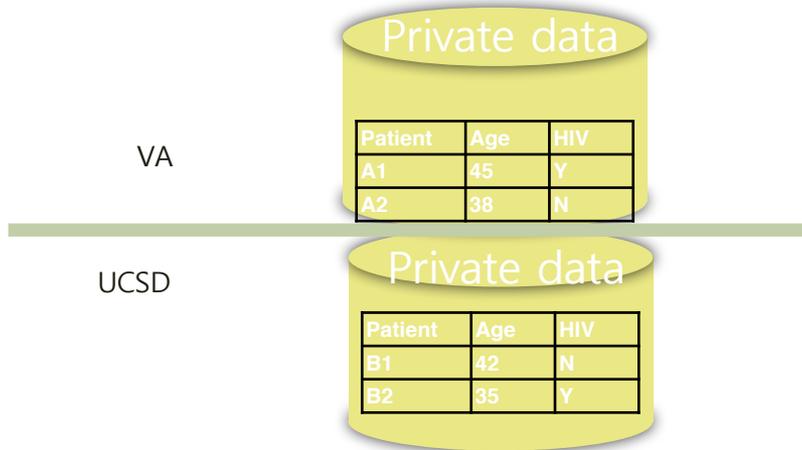




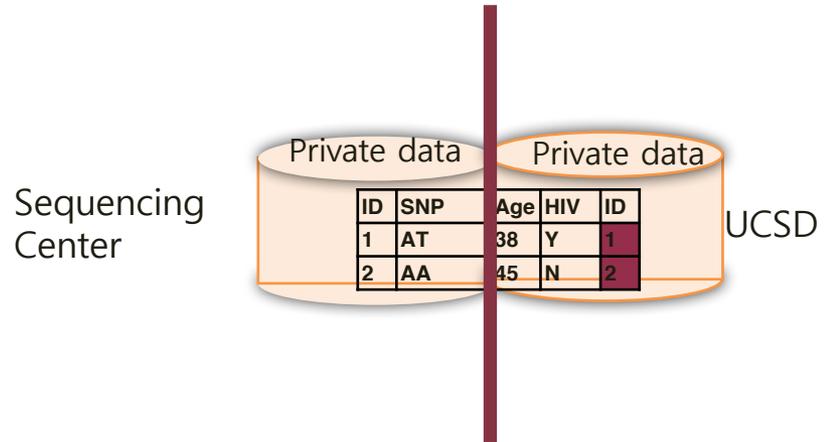
# Privacy-Preserving Distributed Predictive Models

# Two Representative Scenarios

Horizontally distributed data

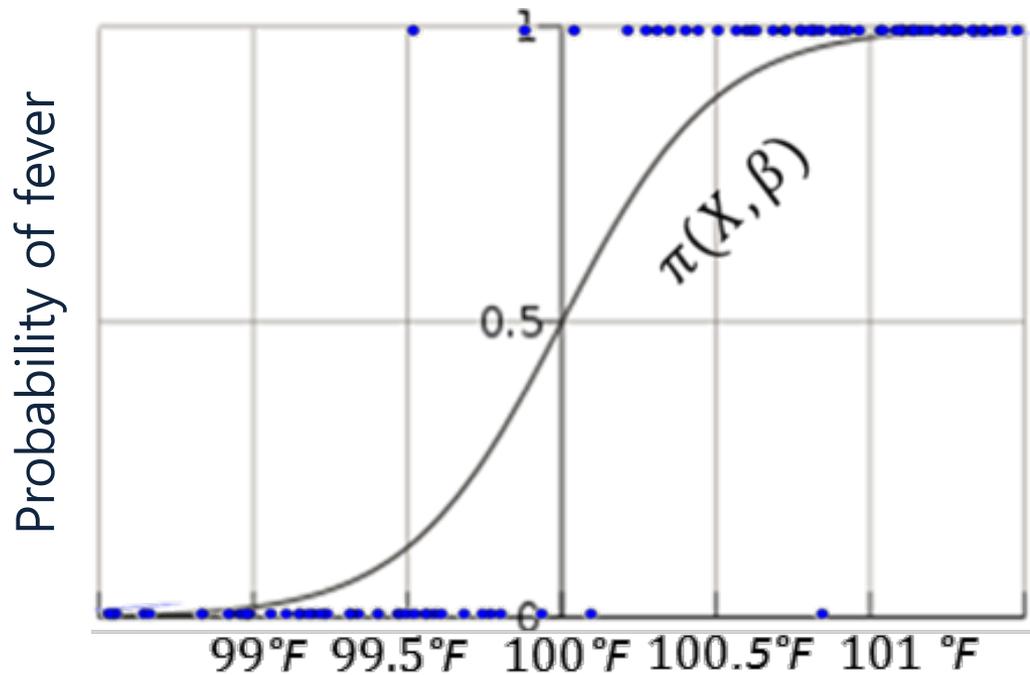


Vertically distributed data



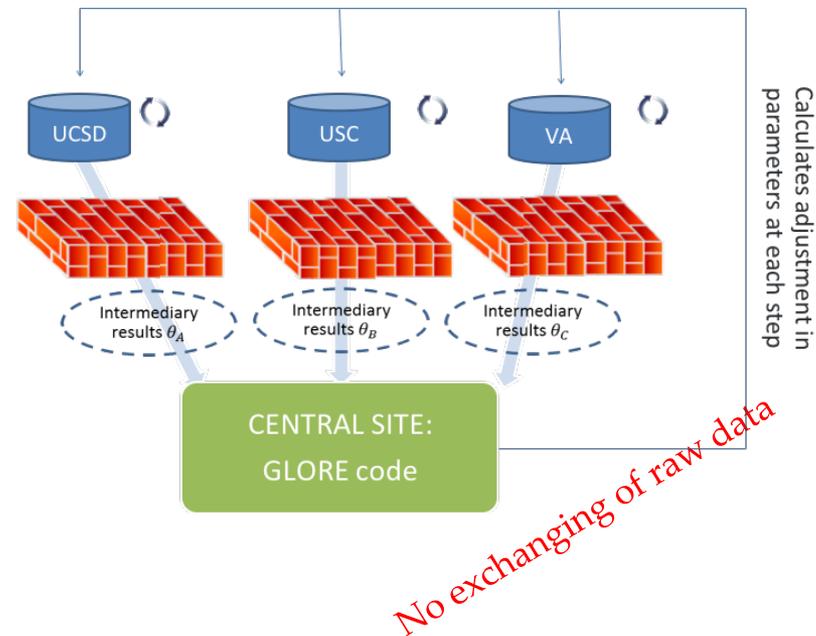
# Logistic Regression

---



# Learning a distributed logistic regression

- Support  $p-1$  features are consistent over  $k$  sites
- In each iteration, intermediary result of a  $p \times p$  matrix and a  $p$ -dimensional vector are transmitted to the central site for optimization



# Maximum Likelihood Estimation

- Estimated probability based on observations of a binary response  $Y$  and covariates  $X$
- Likelihood function based on observed data (centralized)

$$P(Y = 1 | X) = \pi(X, \beta) = \frac{1}{1 + e^{-X\beta}}$$

Binary response      Covariates      Logit function      Model parameter

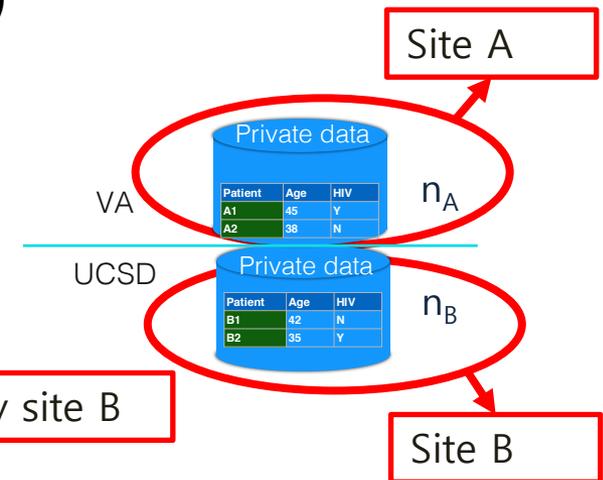
Number of records

$$l(\beta) = \sum_1^n [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

# Maximum Likelihood Estimation

- Likelihood function based on observed data (distributed)

$$P(Y = 1|X) = \pi(X, \beta) = \frac{1}{1 + e^{-X\beta}}$$



Number of records held by site A

$n_A - n_B$

Number of records held by site B

$$l(\beta) = \sum_1^{n_A - n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

# Maximum Likelihood Estimation

---

- Newton-Raphson algorithm for calculation

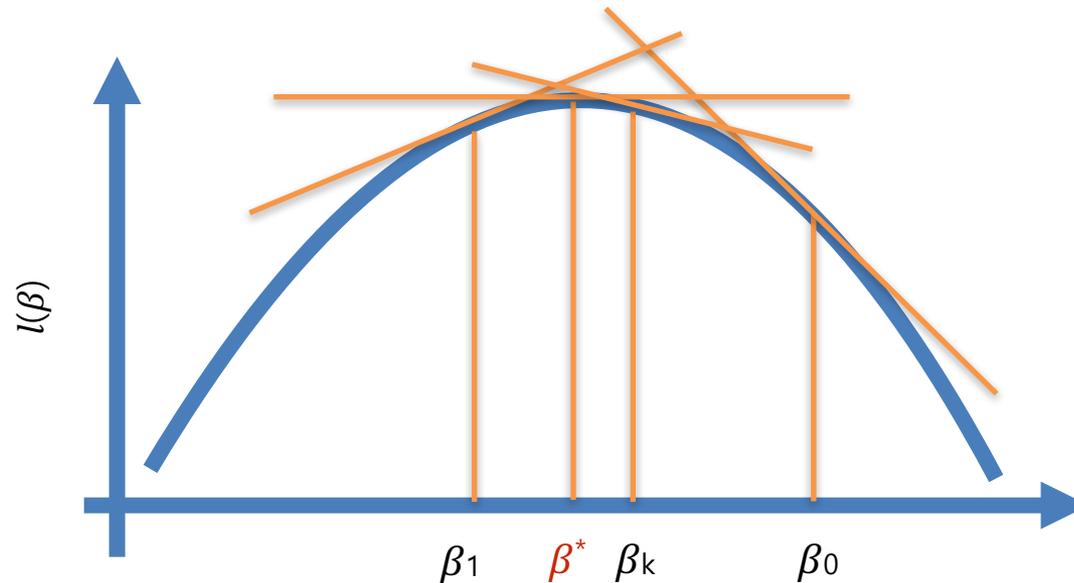
$$P(Y = 1|X) = \pi(X, \beta) = \frac{1}{1 + e^{-X\beta}}$$

$l(\beta)$  is a  
concave  
function

$$\Rightarrow l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

$$\beta^{(k+1)} = \beta^{(k)} - \left[ \frac{\partial^2 l(\beta^{(k)})}{\partial \beta^{(k)} \partial \beta^{(k)T}} \right]^{-1} \frac{\partial l(\beta^{(k)})}{\partial \beta^{(k)}}$$

# Newton-Raphson (NR) Algorithm



$$l(\boldsymbol{\beta}) = \sum_{i=1}^D \{ \boldsymbol{\beta}^T \sum_{l \in \mathcal{D}_i} \mathbf{z}^l - d_i \log[\sum_{l \in \mathcal{R}_i} \exp(\boldsymbol{\beta}^T \mathbf{z}^l)] \}$$

$$\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} - \left[ \frac{\partial^2 l(\boldsymbol{\beta}^{(k)})}{\partial \boldsymbol{\beta}^{(k)} \partial \boldsymbol{\beta}^{(k)T}} \right]^{-1} \frac{\partial l(\boldsymbol{\beta}^{(k)})}{\partial \boldsymbol{\beta}^{(k)}}$$

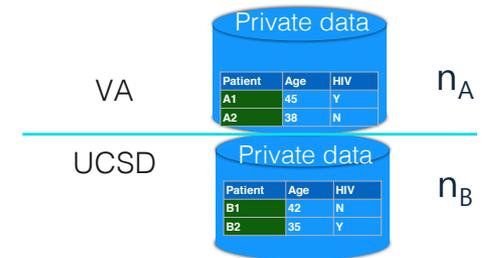
# Distributed Newton-Raphson (NR) Algorithm

---

$$l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

# Distributed Newton-Raphson (NR) Algorithm

$$l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$



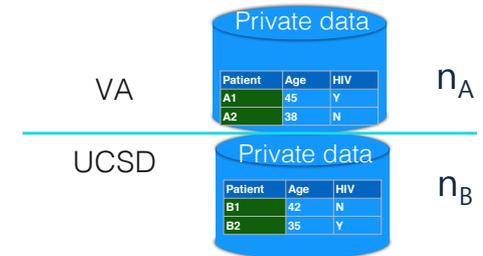
# Distributed Newton-Raphson (NR) Algorithm

$$l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

$$\begin{aligned} \beta^{(k+1)} &= \beta^{(k)} - \left[ \frac{\partial^2 l(\beta^{(k)})}{\partial \beta^{(k)} \partial \beta^{(k)T}} \right]^{-1} \frac{\partial l(\beta^{(k)})}{\partial \beta^{(k)}} \\ &= \beta^{(k)} + [\bar{X}^T \mathcal{W}(\bar{X}, \beta^{(k)}) \bar{X}]^{-1} \bar{X}^T [\bar{Y} - \bar{\Pi}(\bar{X}, \beta^{(k)})] \end{aligned}$$

Global variance-covariance matrix

Global prediction outcomes



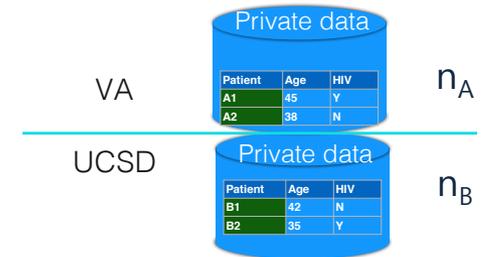
# Distributed Newton-Raphson (NR) Algorithm

$$l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

$$\begin{aligned} \beta^{(k+1)} &= \beta^{(k)} - \left[ \frac{\partial^2 l(\beta^{(k)})}{\partial \beta^{(k)} \partial \beta^{(k)T}} \right]^{-1} \frac{\partial l(\beta^{(k)})}{\partial \beta^{(k)}} \\ &= \beta^{(k)} + [\bar{X}^T W(\bar{X}, \beta^{(k)}) \bar{X}]^{-1} \bar{X}^T [\bar{Y} - \Pi(\bar{X}, \beta^{(k)})] \\ &= \beta^{(k)} + [\bar{X}_A^T W_A(\bar{X}_A, \beta^{(k)}) \bar{X}_A + \bar{X}_B^T W_B(\bar{X}_B, \beta^{(k)}) \bar{X}_B]^{-1} \\ &\quad \cdot \{ \bar{X}_A^T [\bar{Y}_A - \Pi_A(\bar{X}_A, \beta)] + \bar{X}_B^T [\bar{Y}_B - \Pi_B(\bar{X}_B, \beta)] \}. \end{aligned}$$

Local variance-covariance matrix

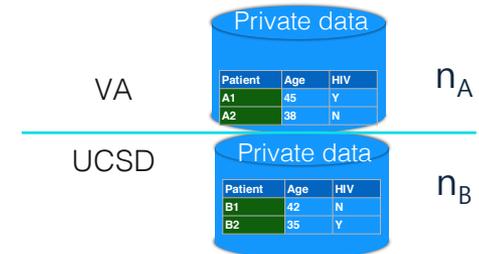
Local prediction outcomes



# Distributed Newton-Raphson (NR) Algorithm

$$l(\beta) = \sum_1^{n_A+n_B} [y_i \log \pi(x_i, \beta) + (1 - y_i) \log(1 - \pi(x_i, \beta))]$$

$$\begin{aligned} \beta^{(k+1)} &= \beta^{(k)} - \left[ \frac{\partial^2 l(\beta^{(k)})}{\partial \beta^{(k)} \partial \beta^{(k)T}} \right]^{-1} \frac{\partial l(\beta^{(k)})}{\partial \beta^{(k)}} \\ &= \beta^{(k)} + [\bar{X}^T W(\bar{X}, \beta^{(k)}) \bar{X}]^{-1} \bar{X}^T [\bar{Y} - \Pi(\bar{X}, \beta^{(k)})] \\ &= \beta^{(k)} + [\bar{X}_A^T W_A(\bar{X}_A, \beta^{(k)}) \bar{X}_A + \bar{X}_B^T W_B(\bar{X}_B, \beta^{(k)}) \bar{X}_B]^{-1} \\ &\quad \cdot \{ \bar{X}_A^T [\bar{Y}_A - \Pi_A(\bar{X}_A, \beta)] + \bar{X}_B^T [\bar{Y}_B - \Pi_B(\bar{X}_B, \beta)] \}. \end{aligned}$$



Local variance-covariance matrix

$$W_A(\bar{X}_A, \beta) = \begin{bmatrix} \pi(x_1, \beta)(1 - \pi(x_1, \beta)) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \pi(x_{n_A}, \beta)(1 - \pi(x_{n_A}, \beta)) \end{bmatrix},$$

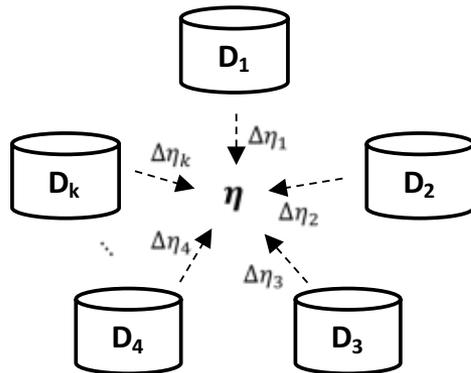
$$W_B(\bar{X}_B, \beta) = \begin{bmatrix} \pi(x_{n_A+1}, \beta)(1 - \pi(x_{n_A+1}, \beta)) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \pi(x_{n_A+n_B}, \beta)(1 - \pi(x_{n_A+n_B}, \beta)) \end{bmatrix},$$

$$\Pi_A(\bar{X}_A, \beta) = \begin{bmatrix} \pi(x_1, \beta) \\ \vdots \\ \pi(x_{n_A}, \beta) \end{bmatrix}, \text{ and } \Pi_B(\bar{X}_B, \beta) = \begin{bmatrix} \pi(x_{n_A+1}, \beta) \\ \vdots \\ \pi(x_{n_A+n_B}, \beta) \end{bmatrix}.$$

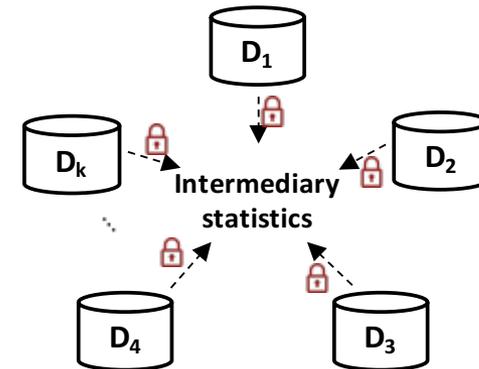
Local prediction outcomes

# What remains to be solved?

- Masking the pattern before transmitting



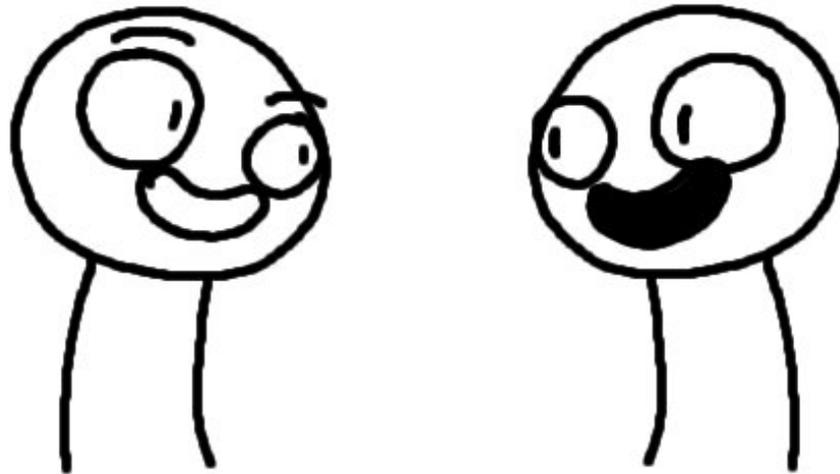
- Using secure primitive to safeguard the communication



# Differential Privacy & homomorphic encryption

---

- A privacy mechanism  $A$  gives  $\epsilon$ -differential privacy if for all neighbouring databases  $D, D'$ , and for any possible output  $S \in \text{Range}(A)$ ,  $\Pr[A(D) = S] \leq \exp(\epsilon) \times \Pr[A(D') = S]$ 
  - $D$  and  $D'$  are **neighboring databases** if they differ on at most one record
- **Homomorphic encryption** is a type of encryption that allows computation conducted on ciphertext, when results are decrypted, map exactly to those of the corresponding computation on the plaintext



# Differential private logistic regression

---

- We perturb the objective function by adding an additional term  $\frac{b^T \beta}{n}$  with  $b$  drawn from a Laplacian distribution with mean 0 and standard deviation  $\frac{2}{\epsilon}$ .

$$\max_{\beta} \left[ l(\beta) = -\sum_{i=1}^n \log(1 + \exp(-y_i \beta^T z_i)) - \frac{\lambda}{2} \beta^T \beta - \frac{b^T \beta}{n} \right]$$

$$\begin{aligned} \beta^{new} &= \beta^{old} - [l''(\beta^{old})]^{-1} l'(\beta^{old}) \\ &= \beta^{old} + (Z^T W^{old} Z + \lambda I)^{-1} \left[ Z^T (Y - \mu^{old}) - \lambda \beta^{old} - \frac{b^T}{n} \right] \end{aligned}$$

$$Z^T W^{old} Z = \sum_k Z_k^T W_k^{old} Z_k, \quad Z^T [Y - \mu^{old}] = \sum_k Z_k^T [Y_k - \mu_k^{old}], \quad b = \sum_k b_k,$$

$$k \in (1, \dots, K)$$

Chaudhuri K, Monteleoni C, Sarwate AD. Differentially Private Empirical Risk Minimization. J Mach Learn Res 2011 Mar;12(Mar):1069–1109. PMID:21892342

# Differentially private logistic regression for distributed data

---

**In a distributed setting, objective perturbation can be achieved by**

- Gamma Distributed Perturbation Laplacian algorithm (DPLA)
- Gauss Distributed Perturbation Laplacian algorithm (DPLA)
- Laplace Distributed Perturbation Laplacian algorithm (DPLA)

Gergely Ács and Claude Castelluccia. I have a DREAM!: differentially private smart metering. In: Proceedings of the 3rd International Conference on Information Hiding. IH'11. 2011, pp. 118-132.  
Goryczka S, Xiong L. A Comprehensive Comparison of Multiparty Secure Additions with Differential Privacy. IEEE Trans Dependable Secure Comput 2017 Sep;14(5):463-477. PMID:28919841

# Differentially private logistic regression for distributed data

---

- **Note that the noise added by a single party is not sufficient to ensure DP!**
  - If we add too much noise, the final output will be less valuable.
  - If we add too little noise, it is not enough to protect the privacy.
- **Privacy mechanisms are not designed to provide security of computations.**
  - We need to protect the intermediary results, otherwise, privacy cannot be ensured in a global manner

---

# Algorithm

---

# Win-Win Strategy

$$\beta^{new} = \beta^{old} - [l''(\beta^{old})]^{-1} l'(\beta^{old})$$

$$= \beta^{old} + \underbrace{\left( \sum_k Z_k^T W_k^{old} Z_k + \frac{\lambda}{K} I \right)^{-1}}_{\text{Hessian} = H} \underbrace{\left[ \sum_k \left( Z_k^T [Y_k - \mu_k^{old}] - \frac{\lambda}{K} \beta^{old} \right) - \left( \sum_k \frac{1}{n} b_k \right)^T \right]}_{\text{Gradient} = g}$$

- **Homomorphic Encryption with “Fixed Hessian”**

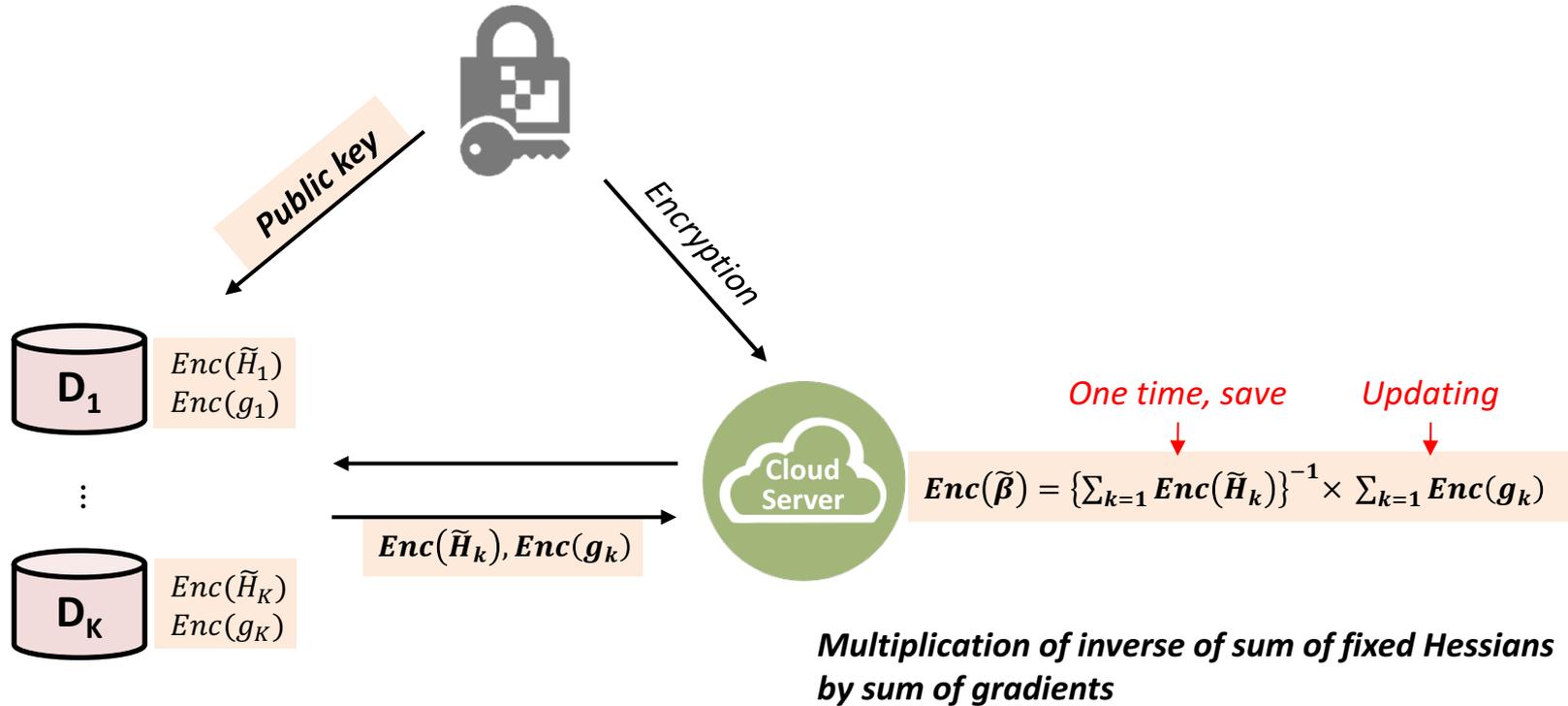
- $\sum_k Z_k^T W_k^{old} Z_k + \frac{\lambda}{K} I = \sum_k H_k \approx \sum_k \frac{1}{4} Z_k^T Z_k + \frac{\lambda}{K} I = \sum_k \bar{H}_k \approx \sum_k \text{diag}(\bar{H}_k) = \sum_k \tilde{H}_k$  *One time*  $\text{Enc}(\tilde{H}_k)$  
- $Z^T [Y - \mu^{old}] - \lambda \beta^{old} - \frac{1}{n} b = \sum_k Z_k^T [Y_k - \mu_k^{old}] - \frac{\lambda}{K} \beta^{old} - \sum_k \frac{1}{n} b_k = \sum_k g_k$  *Iteratively*  $\text{Enc}(g_k)$  

- **Differential Privacy**

- $\beta^{new}$  can be revealed to parties because of the noise
- HE can be renewed every iteration

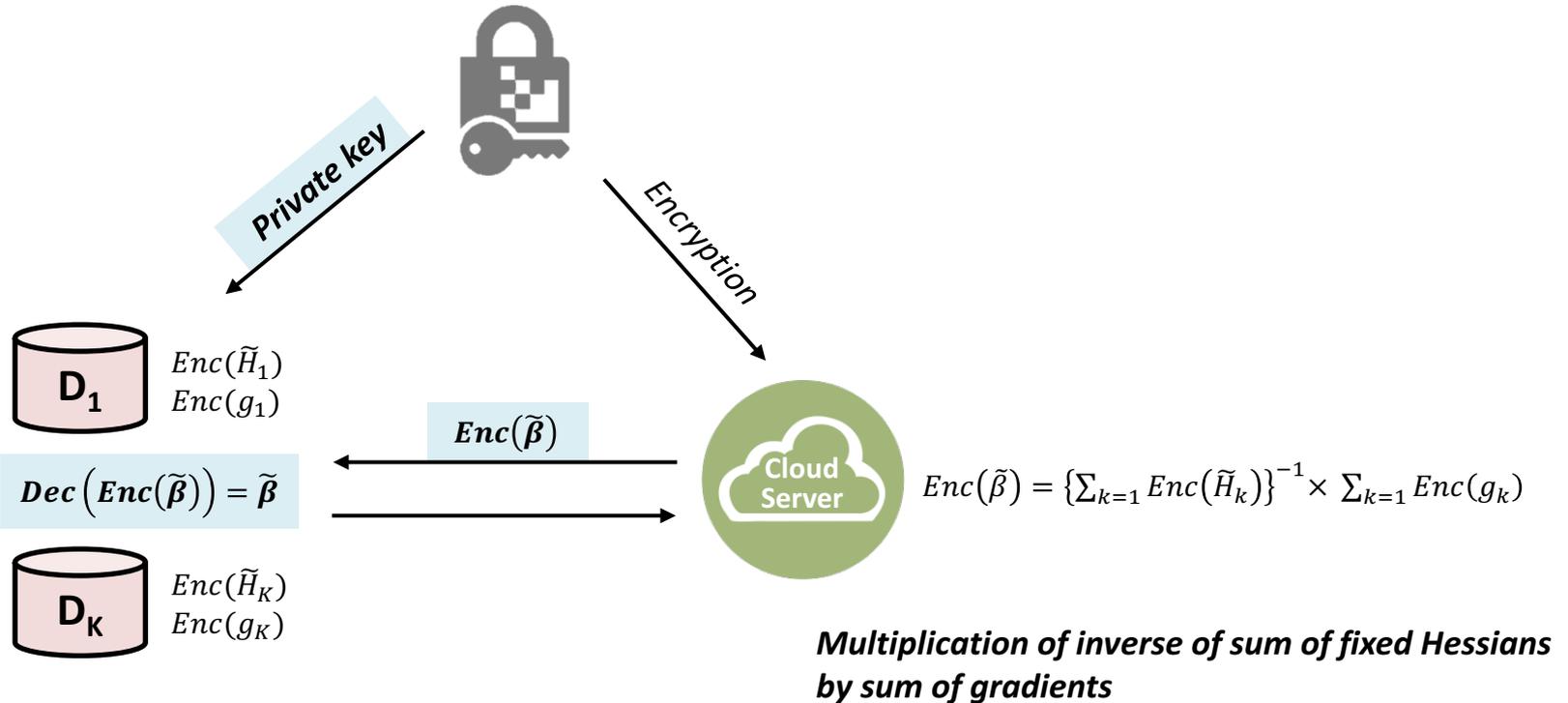
Based on DP, we can reduce time complexity and error accumulation of HE

# SMC Schemes with HE under Fixed Hessian



*Approximation of fixed Hessian and gradient*

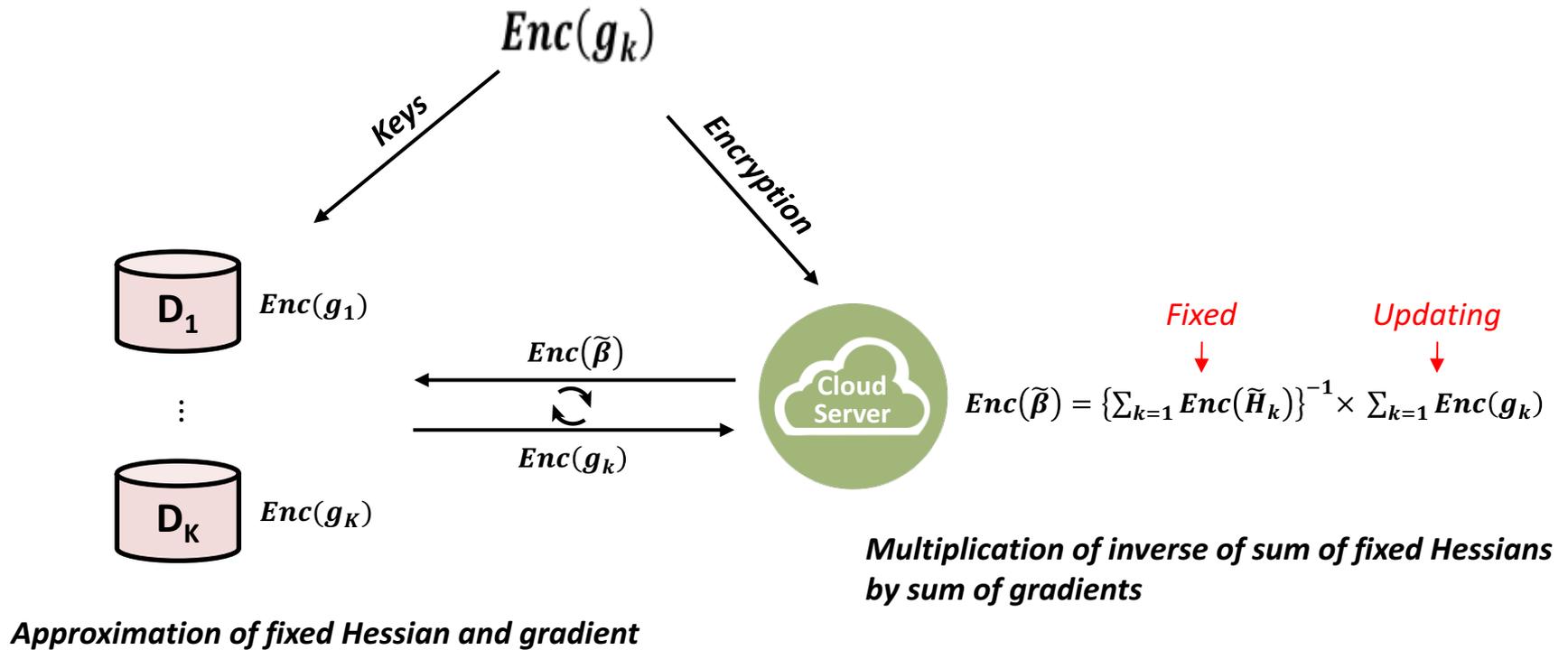
# SMC Schemes with HE under Fixed Hessian



*Approximation of fixed Hessian and gradient*

$\tilde{\beta}$  can be revealed to parties because of the noise.

# SMC Schemes with HE under Fixed Hessian



A few iterations → Converge

# Limitations of Fixed Hessian

---

$$\sum_k Z_k^T W_k^{old} Z_k + \frac{\lambda}{K} I = \sum_k H_k \approx \sum_k \frac{1}{4} Z_k^T Z_k + \frac{\lambda}{K} I = \sum_k \bar{H}_k \approx \sum_k \text{diag}(\bar{H}_k) = \sum_k \tilde{H}_k$$

- Simple approximation of Hessian using only its diagonal elements
- Valid when the matrix strongly diagonally dominant
- Large enough  $\lambda$  to be set

Largely dependent on  $\lambda$

- Better diagonal Hessian approximation

Diagonal Updating via Quasi-Cauchy Relation

# Diagonal Updating via Quasi-Cauchy Relation\*

---

$$\nabla^2 f(x) = \nabla^2 f_A(x) + \nabla^2 f_B(x)$$

where  $\nabla^2 f_A(x)$ : a diagonal matrix consisting the diagonal entries of the Hessian

$\nabla^2 f_B(x)$ : the actual Hessian except that its diagonal entries are all zero

$$\nabla^2 f(x) \approx D = \Psi_1 + \Psi_2 = \Psi_1 + (\theta I + \Psi_3)$$

where  $\Psi_1$ : a positive definite diagonal matrix

$$\min \frac{1}{2} \|\Psi_3\|_F^2,$$

$$s.t. \quad s_i^T (\Psi_1 + (\theta I + \Psi_3)) s_i = s_i^T y_i \text{ and } \Psi_3 \text{ is diagonal}$$

$$D_{i+1} = D_i + \frac{s_i^T y_i - s_i^T \Psi_1 s_i - \theta s_i^T s_i}{\text{tr}(E_i^2)} E_i$$

where  $\theta_i = \min \left[ 1, \frac{s_i^T y_i - s_i^T \Psi_1 s_i}{s_i^T s_i} \right]$  for positive definiteness and  $E_i = \text{diag}(s_{i,1}^2, s_{i,2}^2, \dots, s_{i,m}^2)$

\* Marjugi and Leong (2013) Diagonal Hessian Approximation for Limited Memory Quasi-Newton via Variational Principle, *Journal of Applied Mathematics*

# Diagonal Updating via Quasi-Cauchy Relation

*Decomposable*

$$D_{i+1} = D_i + \frac{s_i^T y_i - s_i^T \Psi_1 s_i - \theta_i s_i^T s_i}{\text{tr}(E_i^2)} E_i, \quad \boxed{\frac{s_i^T y_i - s_i^T \Psi_1 s_i}{\text{tr}(E_i^2)} E_i} - \theta_i \cdot \frac{s_i^T s_i}{\text{tr}(E_i^2)} E_i = \sum_{k=1}^K V_{ik} - \theta_i \cdot W_i$$

where  $s_i = \beta^{i+1} - \beta^i$

$$y_i = \sum_k \left( Z_k^T [Y_k - \mu_k^{i+1}] - \frac{\lambda}{K} \beta^{i+1} \right) - \sum_k \left( Z_k^T [Y_k - \mu_k^i] - \frac{\lambda}{K} \beta^i \right)$$

$$V_{ik} = \frac{s_{ik}^T y_i - s_{ik}^T \Psi_1 s_{ik}}{\text{tr}(E_i^2)} E_i, \quad W_i = \frac{s_i^T s_i}{\text{tr}(E_i^2)} E_i$$

- For positive definiteness
  - $\theta_i = \min \left[ 1, \frac{s_i^T y_i - s_i^T \Psi_1 s_i}{s_i^T s_i} \right]$
  - Comparison within ciphertext is not easy, so we used one more round of iteration

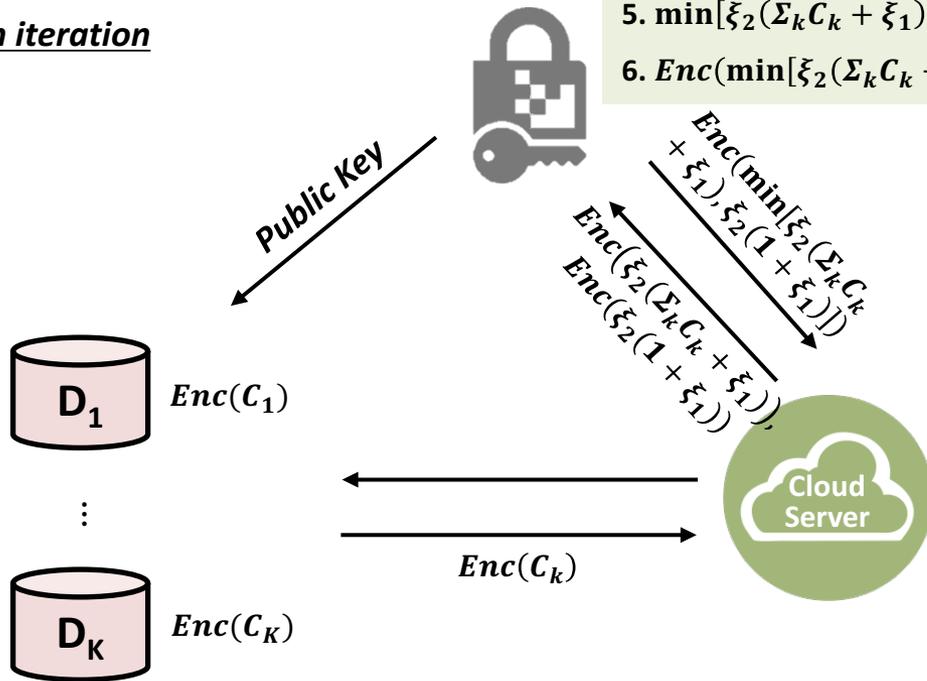
**One more step is added every iteration.**

# Positive Definiteness

$$\theta_i = \min \left[ 1, \frac{\text{Decomposable} \left[ \frac{s_i^T y_i - s_i^T \Psi_1 s_i}{s_i^T s_i} \right]}{s_i^T s_i} \right] = \min[1, \sum_k C_{ik}]$$

$$\text{where } C_{ik} = \frac{s_i^T y_{ik}}{s_i^T s_i} - \frac{s_i^T \Psi_1 s_i}{K s_i^T s_i}$$

*i*-th iteration



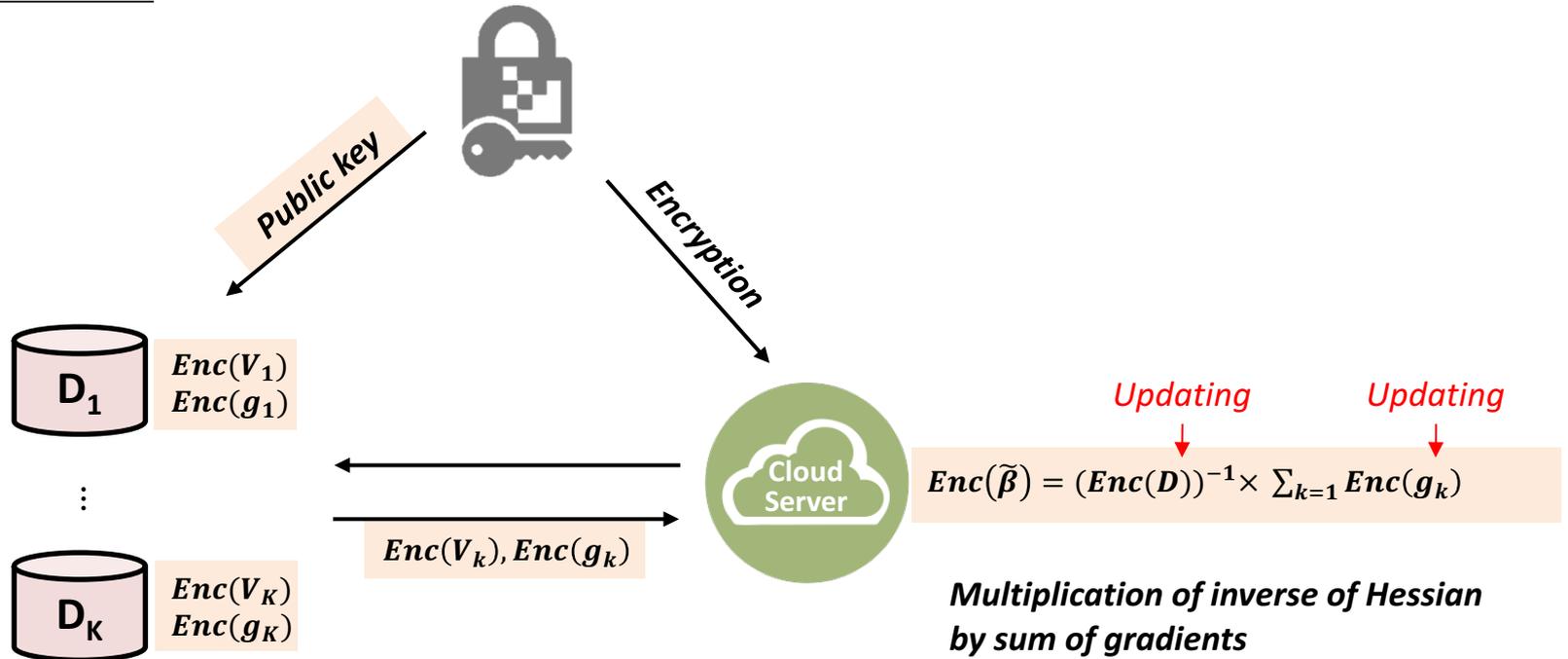
4.  $Dec(\xi_2(\sum_k C_k + \xi_1)) = \xi_2(\sum_k C_k + \xi_1)$   
 $Dec(\xi_2(1 + \xi_1)) = \xi_2(1 + \xi_1)$
5.  $\min[\xi_2(\sum_k C_k + \xi_1), \xi_2(1 + \xi_1)]$
6.  $Enc(\min[\xi_2(\sum_k C_k + \xi_1), \xi_2(1 + \xi_1)])$

1.  $Enc(\sum_k C_k) = \sum_k Enc(C_k)$
2. Random number  $\xi_1, \xi_2$  generation
3.  $Enc(\xi_2(\sum_k C_k + \xi_1)), Enc(\xi_2(1 + \xi_1))$

7.  $Enc(\min[\sum_k C_k, 1])$   
 $= Enc(\min[\xi_2(\sum_k C_k - \xi_1), \xi_2(1 + \xi_1)])$   
 $/ Enc(\xi_2^{-1}) - Enc(\xi_1)$

# SMC Schemes with HE under Updating Hessian

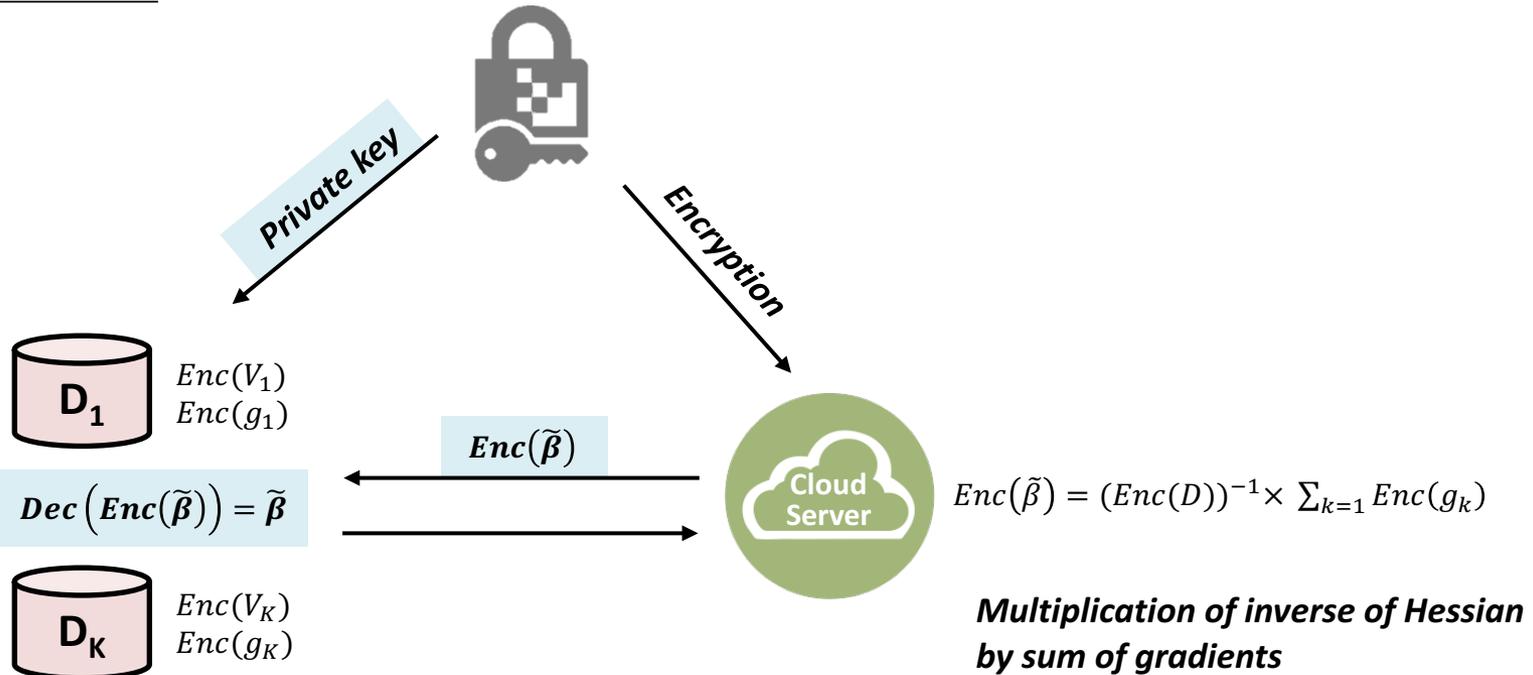
*i*-th iteration



*Approximation of Hessian and gradient*

# SMC Schemes with HE under Updating Hessian

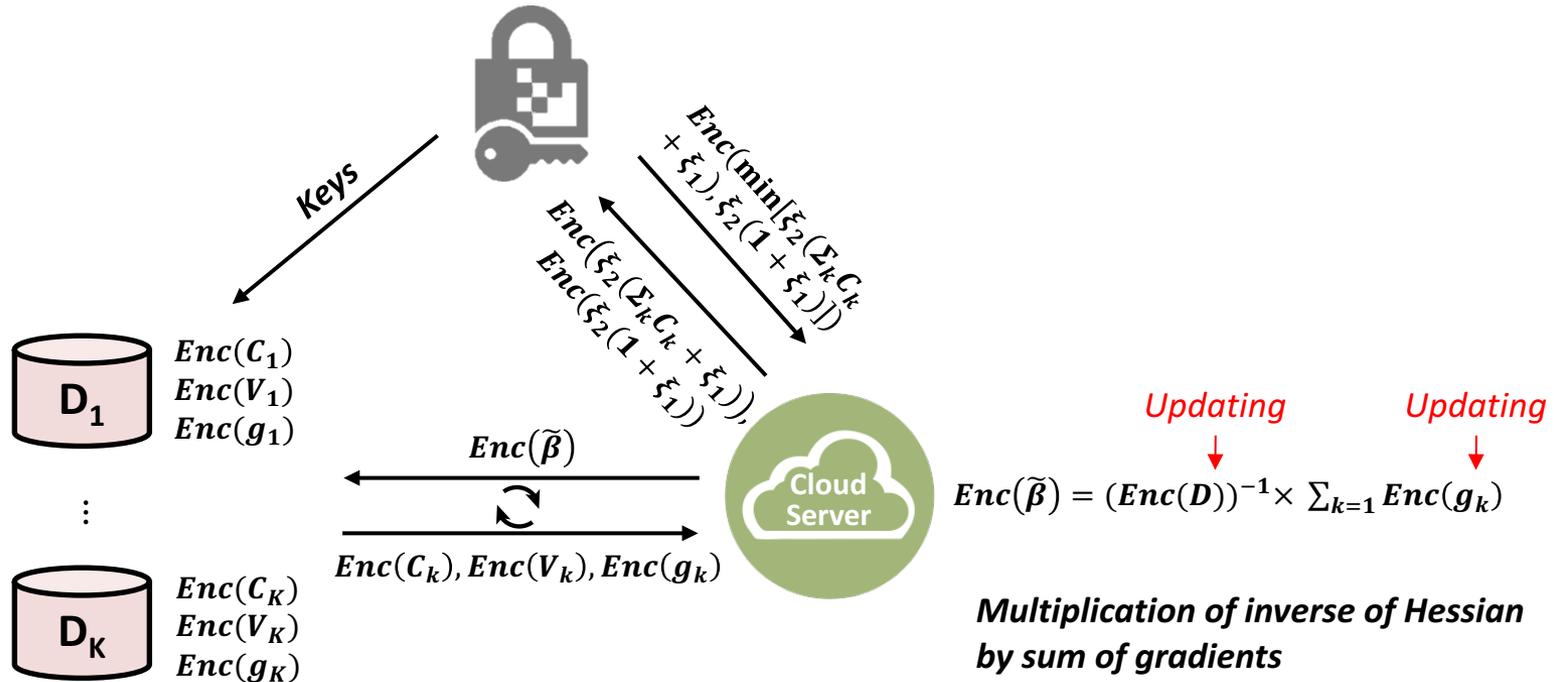
*i*-th iteration



*Approximation of Hessian and gradient*

$\tilde{\beta}$  can be revealed to parties because of the noise.

# SMC Schemes with HE under Updating Hessian



*Approximation of Hessian and gradient*

**A few iterations → Converge**

# Trade-Off between Fixed Hessian and Updating Hessian

---

	Fixed Hessian	Updating Hessian
The number of Iterations	↑	↓
Time per iteration	↓	↑

- **Fixed Hessian**

- Iteration numbers are too dependent on  $\lambda$  which is also depending on data.
- The number of iterations can be more than 100 when not big enough  $\lambda$ .

- **Updating Hessian**

- Iteration numbers are quite robust on  $\lambda$ .
- Inverse diagonal Hessian is not that much expensive ( $\because$  vector computation).

# Dataset: Death in hospital

---

- **PhysioNet Challenge 2012 [MIMIC II database]**

- Dataset comprised of 4000 patient stays in the ICU lasting at least 2 days for predicting mortality.
  - The data were formatted as time-stamped measurements for 37 distinct variables.
  - Four static variables (age, gender, height, and initial weight) are also present.
- Number of patients: 4000, Number of features: 41

# Data Preprocessing

Percentage of patients for whom at least one measurement was available during the first 48 ICU hours

Measurement	%	Measurement	%
ABP (Arterial blood pressure)		Heart rate	98.4
Invasive (diastolic, mean, systolic)	98.4	K (Serum potassium)	97.9
Non-invasive (diastolic)	87.3	Lactate	54.8
Non-invasive (mean)	87.2	Mg (Serum magnesium)	97.5
Non-invasive (systolic)	87.6	Mechanical ventilation	63.1
Albumin	40.5	Na (Serum sodium)	98.2
ALP (Alkaline phosphatase)	42.4	PaCO2	75.4
ALT (Alkaline transaminase)	43.4	PaO2	75.4
AST (Aspartate transaminase)	43.4	pH	75.9
Bilirubin	43.4	Platelets	98.3
BUN (Blood urea nitrogen)	98.4	Respiration rate	27.7
Cholesterol	7.9	SaO2	44.7
Creatinine	98.4	Temperature	98.4
FiO2 (Fractional inspired oxygen)	67.6	Troponin-I	4.7
Glasgow Coma Score (GCS)	98.4	Troponin-T	21.9
Glucose	97.5	Urine output	97.4
HCO3 (Serum bicarbonate)	98.2	WBC (White blood cell count)	98.2
HCT (Hematocrit)	98.4	Weight	67.7

1. Compute min, max, mean, first value, last value as a way to represent time-series features
2. Missing values are replaced by the mean value of a feature.

**Number of patients: 4,000, Number of features: 189**

# Experiment

---

- **Models**

- 1) (Distributed) model without differential privacy

For  $b \sim Lap(0, \sqrt{2}/\epsilon)$  with  $2/\epsilon$  standard deviation

- 2) Model with Gamma DLPA and HE

- 3) Model with Gauss DLPA and HE

- 4) Model with Laplace DLPA and HE

- **Scenario**

- 3 sites with equal sizes

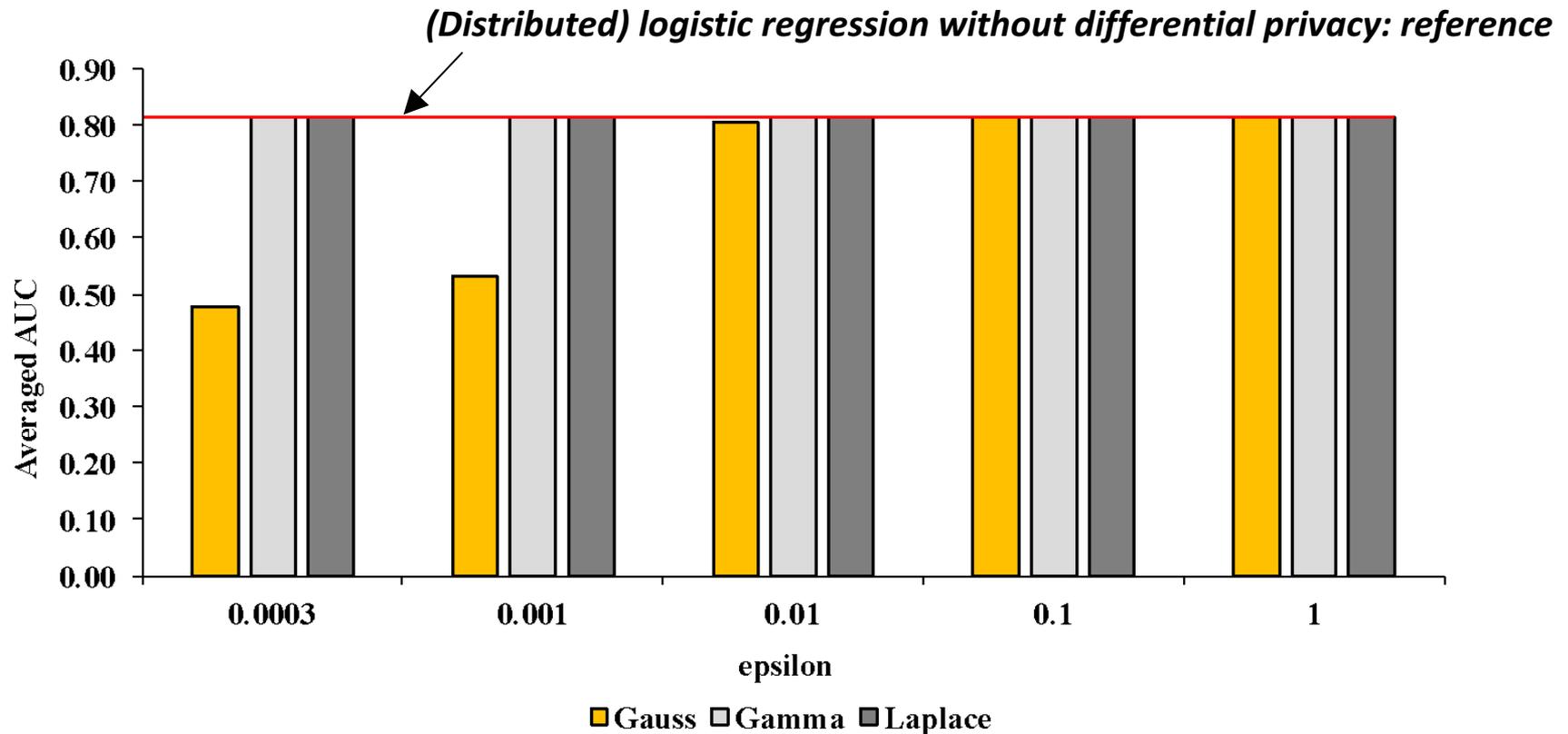
- **Comparison**

- 10 repetitions of 4-fold CV

- AUC
- Mean of coefficients
- Standard deviation of coefficients

# Reference Result without HE

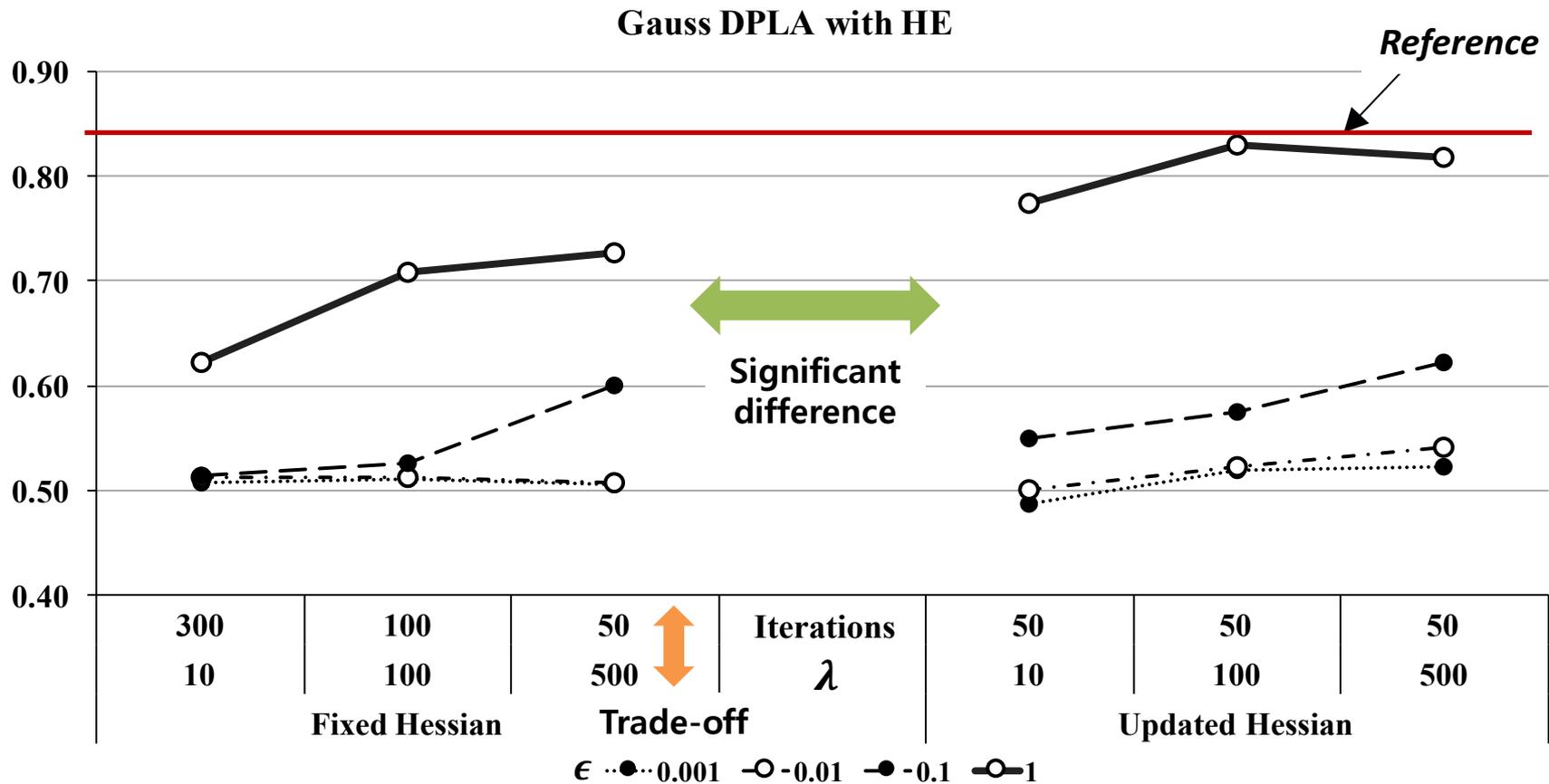
- Averaged AUC on plaintext



Gauss  $\ll$  Laplace  $\approx$  Gamma

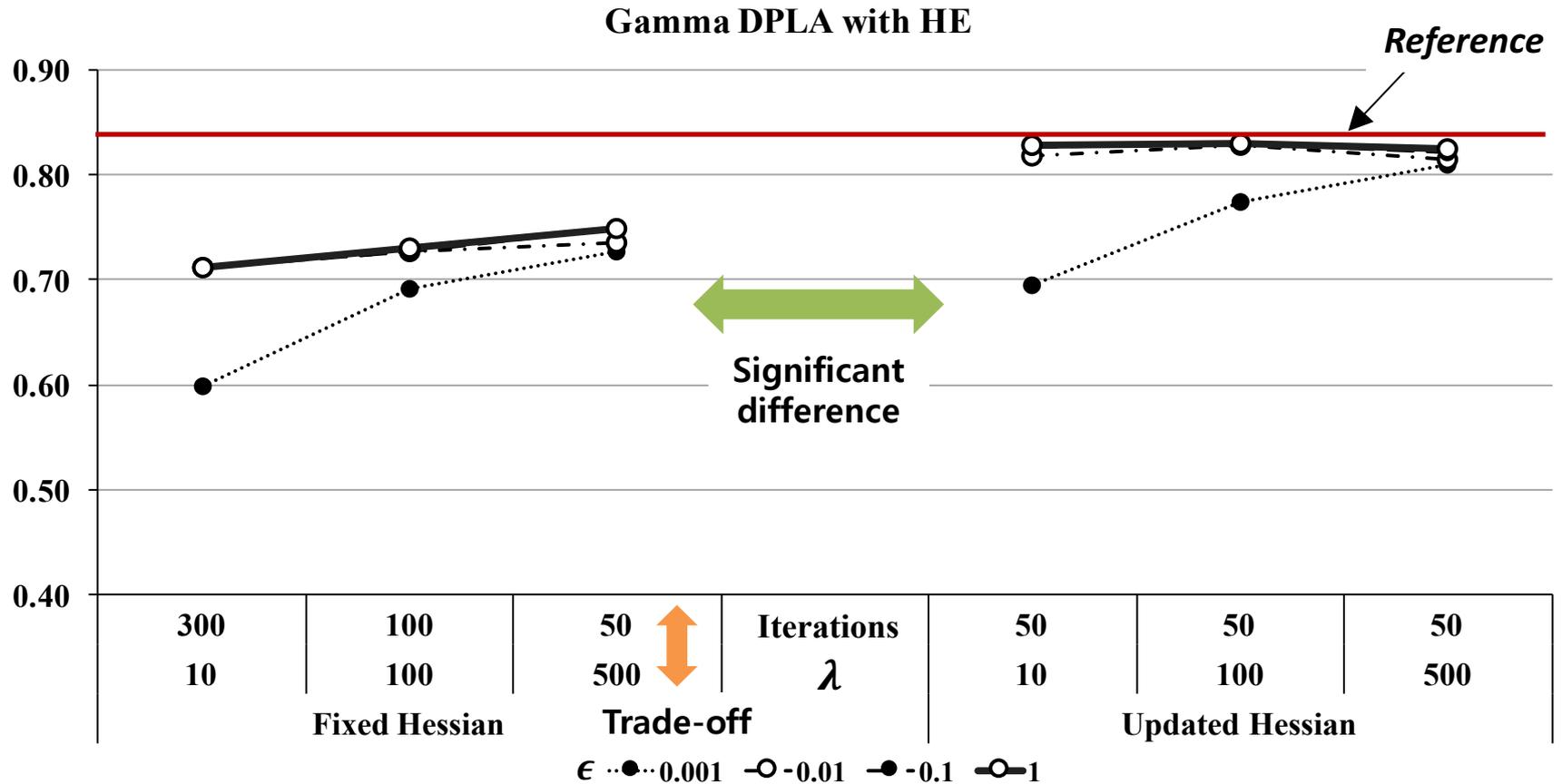
# Gauss DPLA with HE

Budget:  $\epsilon$ /iterations



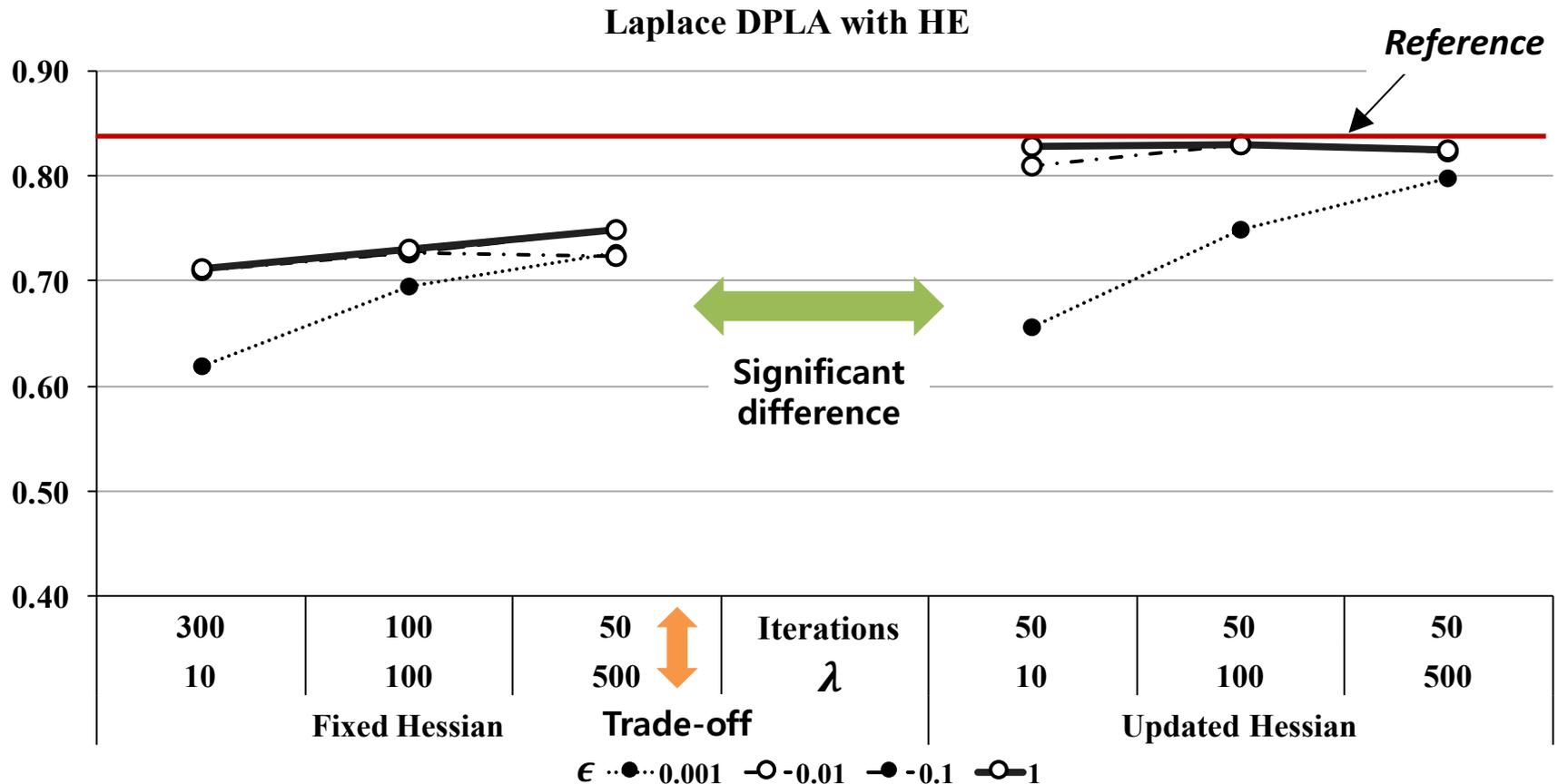
# Gamma DPLA with HE

Budget:  $\epsilon$ /iterations



# Laplace DPLA with HE

Budget:  $\epsilon$ /iterations



Same tendency: Gauss  $\ll$  Laplace  $\approx$  Gamma

# Time Complexity

	Fixed Hessian	Updating Hessian
The number of Iterations	↑	↓
Time per iteration	↓	↑

	$\lambda$		
	10	100	500
	<b>Fixed Hessian</b>		
Iterations	50	100	300
Time (s)	<b>26.33</b>	49.36	<b>142.39</b>
Trade-off depending on the number of iterations			
	<b>Updated Hessian</b>		
Iterations	50	50	50
Time (s)	<b>131.72</b>	<b>132.23</b>	<b>131.22</b>

**We confirmed win-win strategy!**

# Acknowledgement

---

- Junghye Lee
  - Miran Kim
  - Shuang Wang
  - Robert El-Kareh
  - Lucila Ohno-Machado
- NIGMS grant  
R01GM118609

# Q&A

---

