

# 遺伝情報と医療情報のアウトソーシング型 プライバシ保護統計解析

佐久間 淳

筑波大学/JST CREST/理研AIP

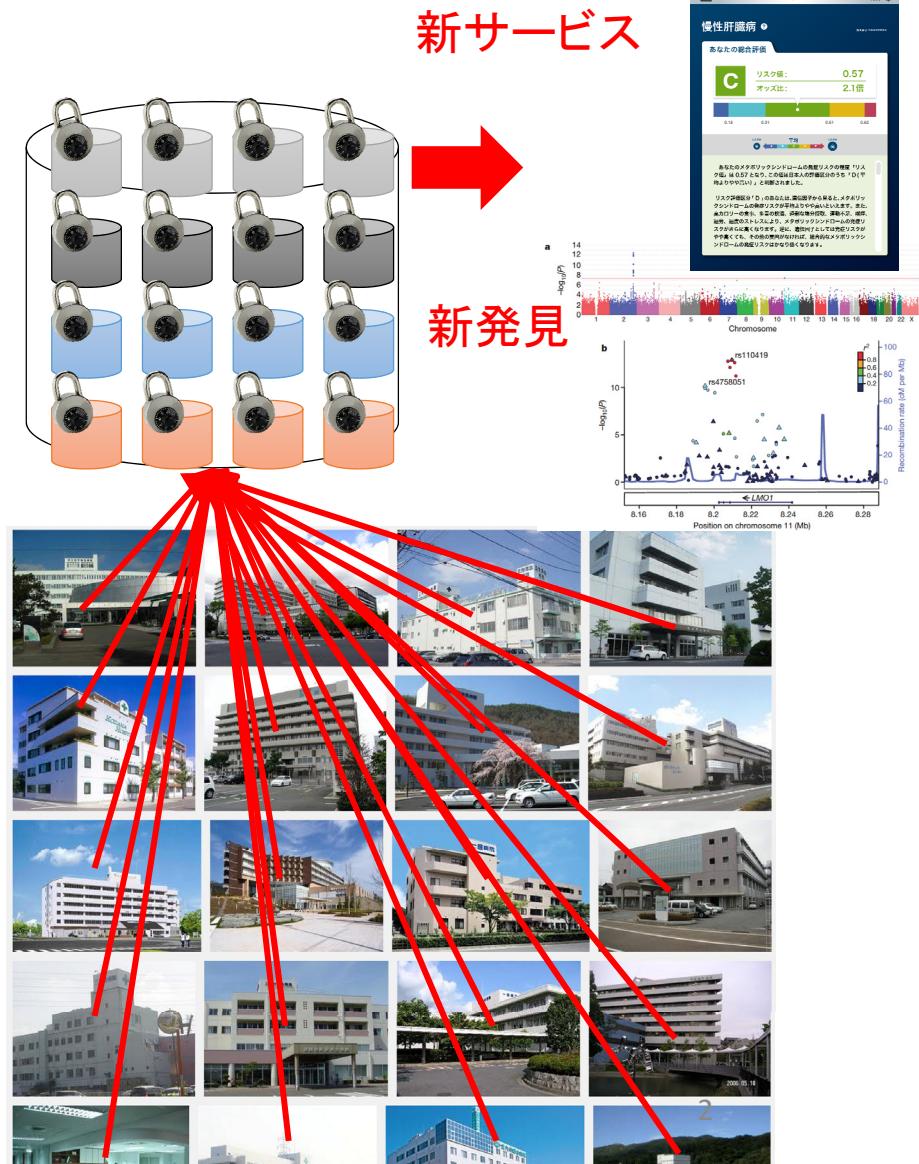


Joint work with

陸文傑, 川崎将平, 中澤貴明, 西出隆志(筑波)  
國廣昇, 津田宏治、竹内聖悟(東大)  
竹内一郎、松井孝太、花田博幸、高田敏行(名工大)  
山田芳司、安河内 彦輝(三重大)

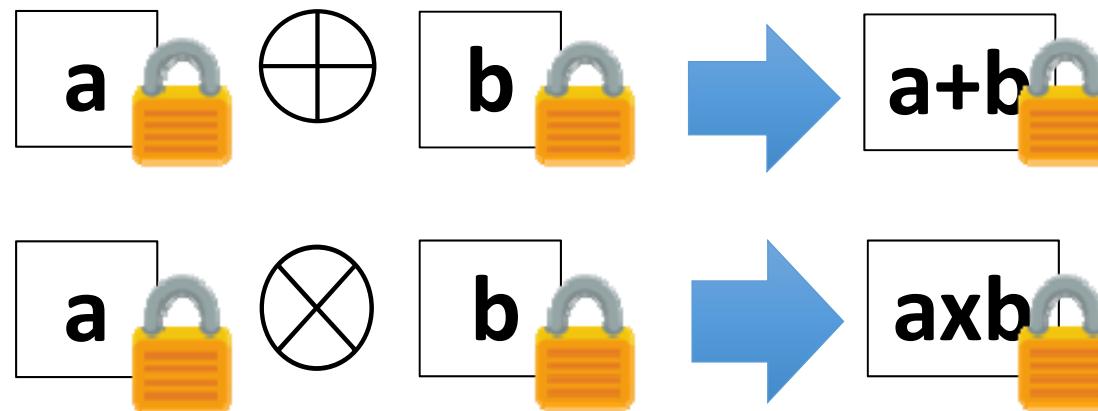
# 目標:孤立small dataをBIG DATAに

孤立したSmall data



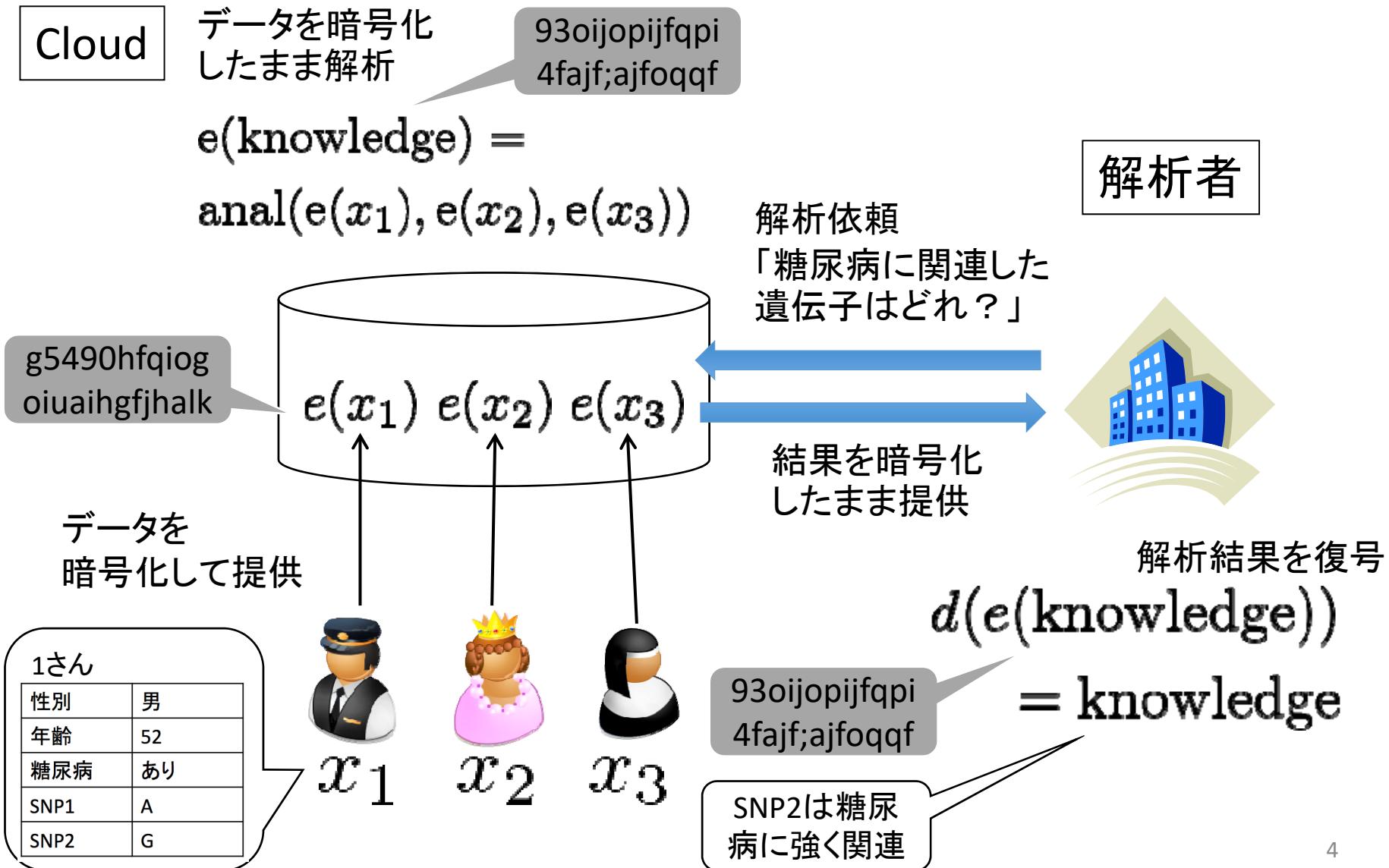
# 完全準同型暗号

- ・データを暗号化したまま加算・乗算
- ・原理的には任意の演算を「データを見せずに」評価可能

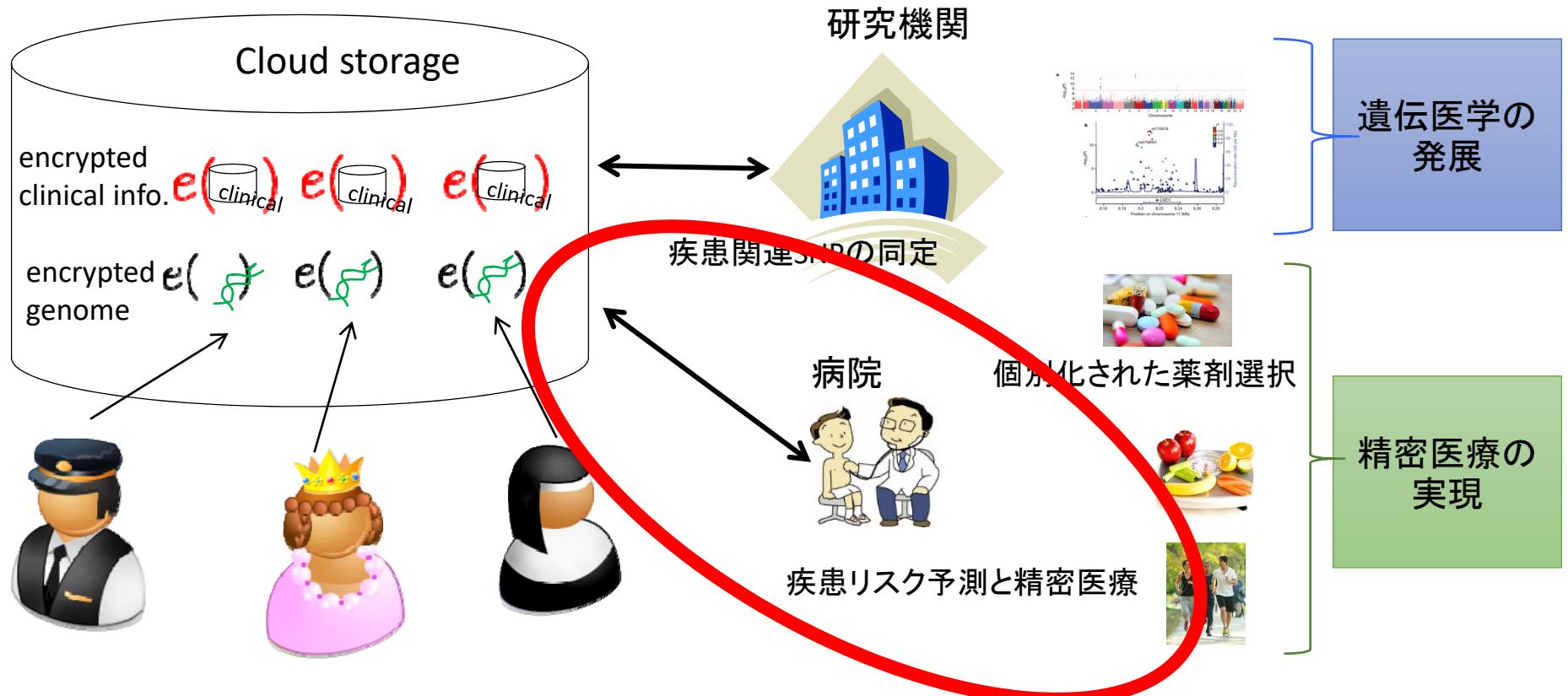


Gentry Craig, "A fully homomorphic encryption scheme",  
Doctoral dissertation, Stanford University, 2009

# 基盤技術：秘密計算・準同型暗号

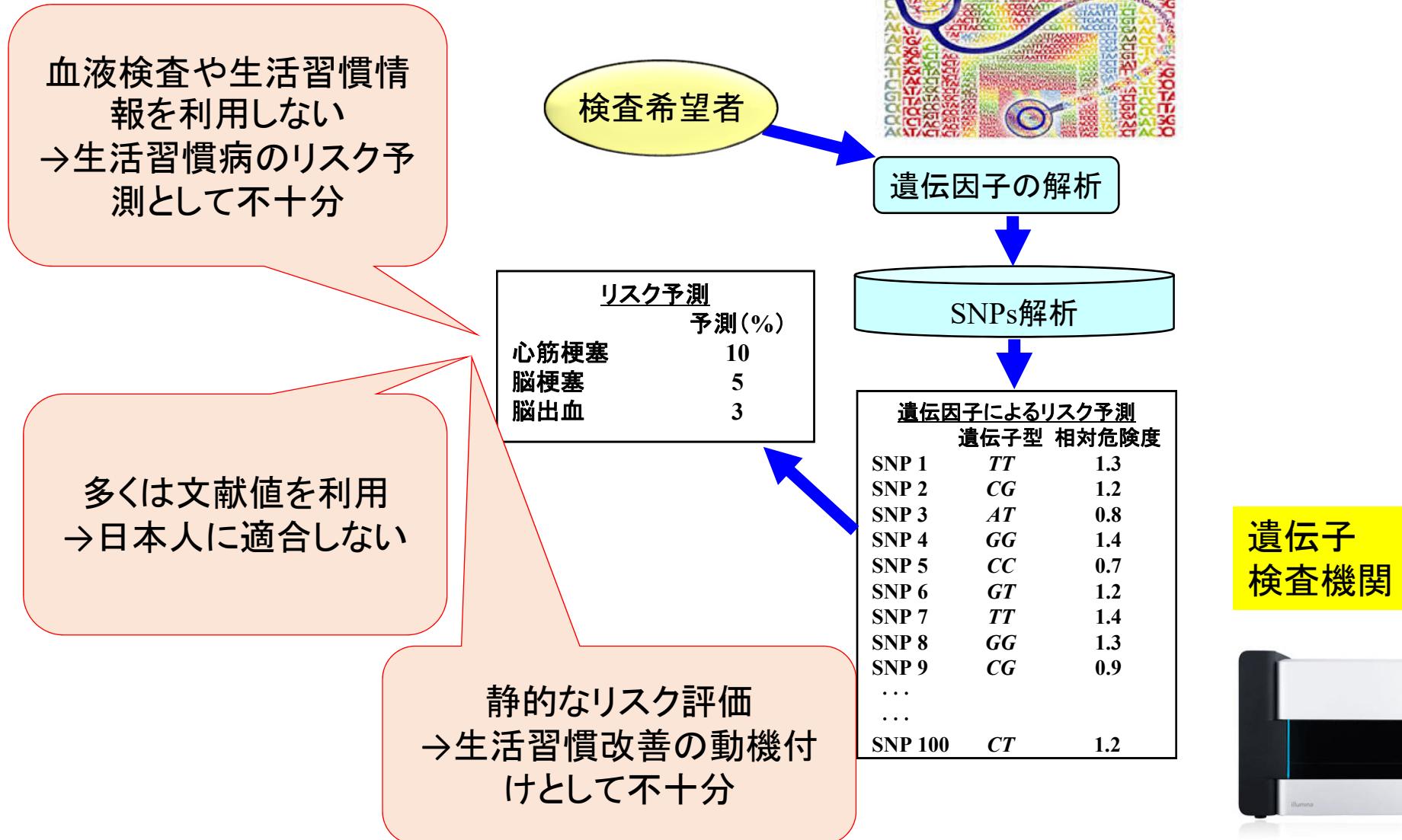


# CRESTにおける研究開発: 暗号化個人ゲノム利用基盤

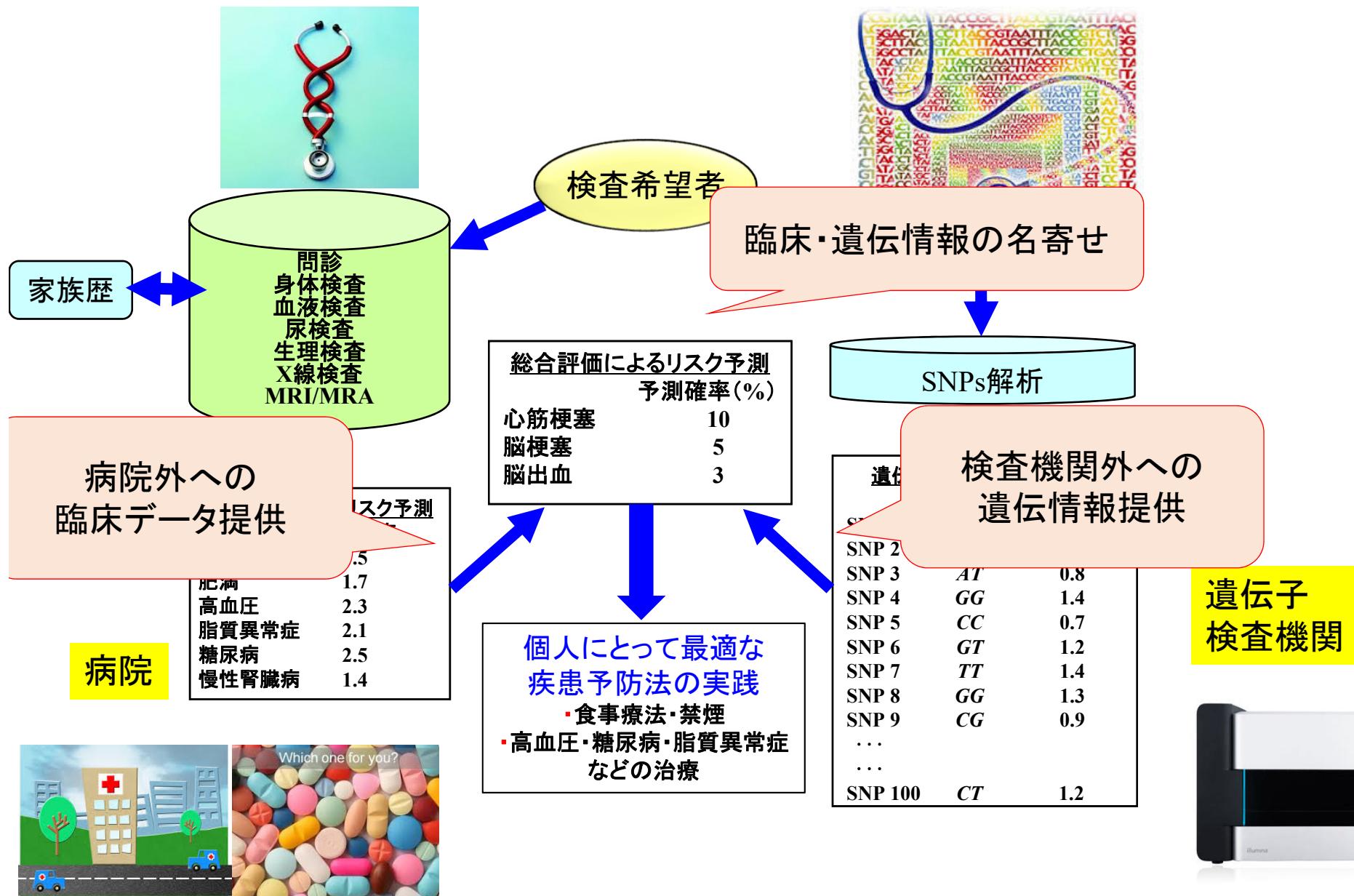


- ・個人ゲノムは解析は高コストだが生涯変化しない
- ・準同型暗号化個人ゲノムをクラウド上に蓄積
- ・様々な機関から様々な用途で安全に利用できるように個人ゲノム解析のサービス提供

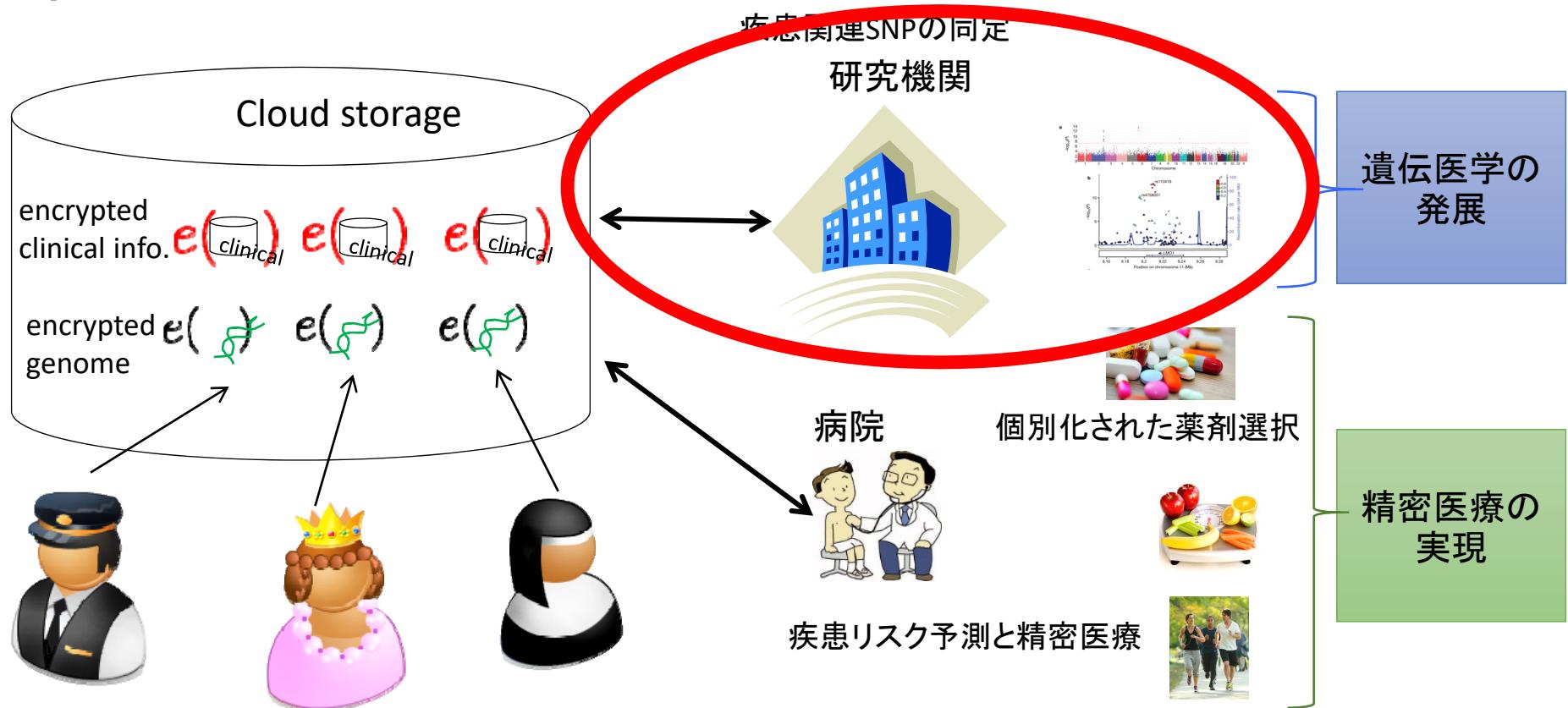
# 従来の生活習慣病個別化予防システム



# 日本人のための生活習慣病個別化予防システム



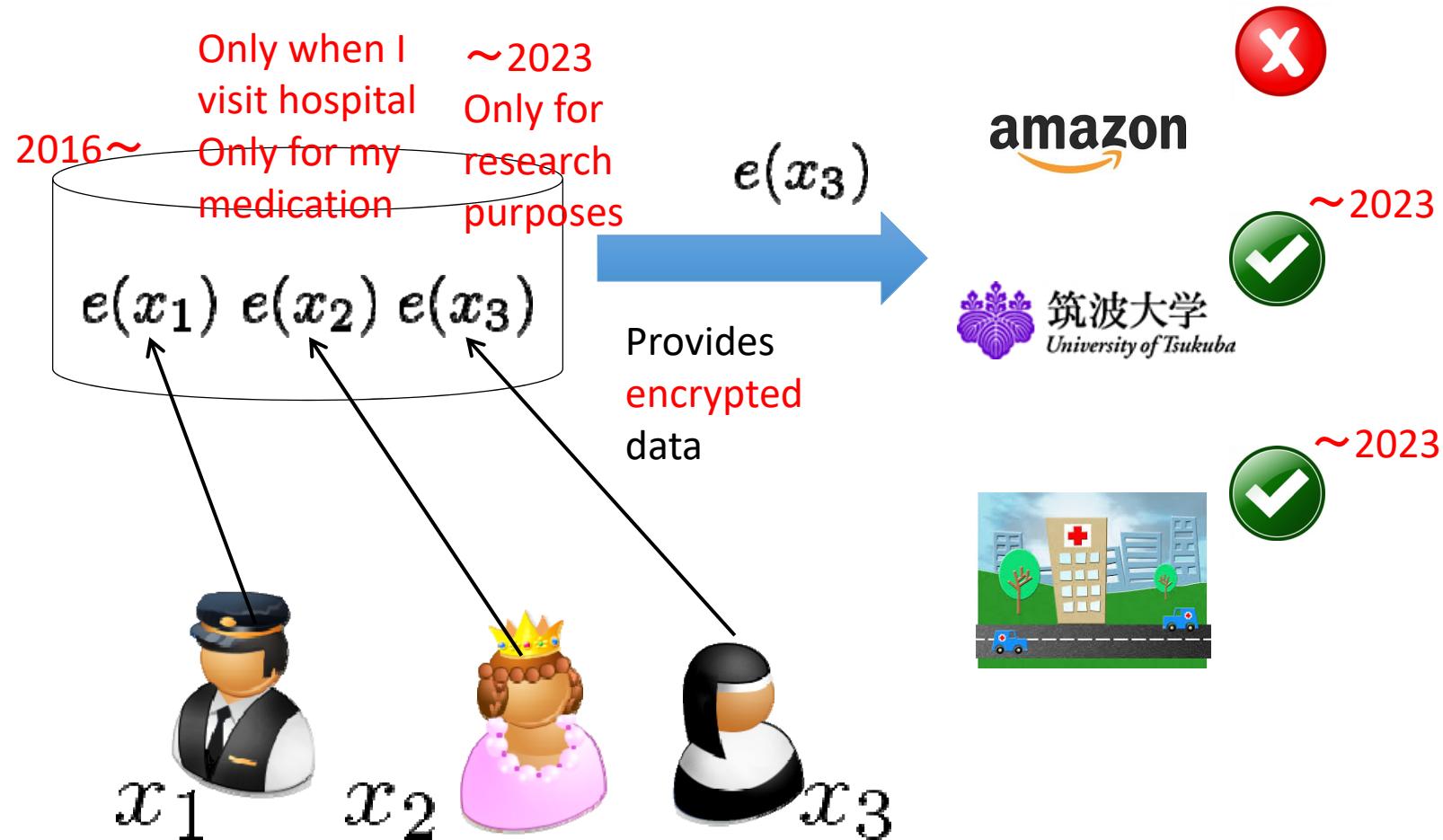
# 暗号化個人ゲノム利用基盤



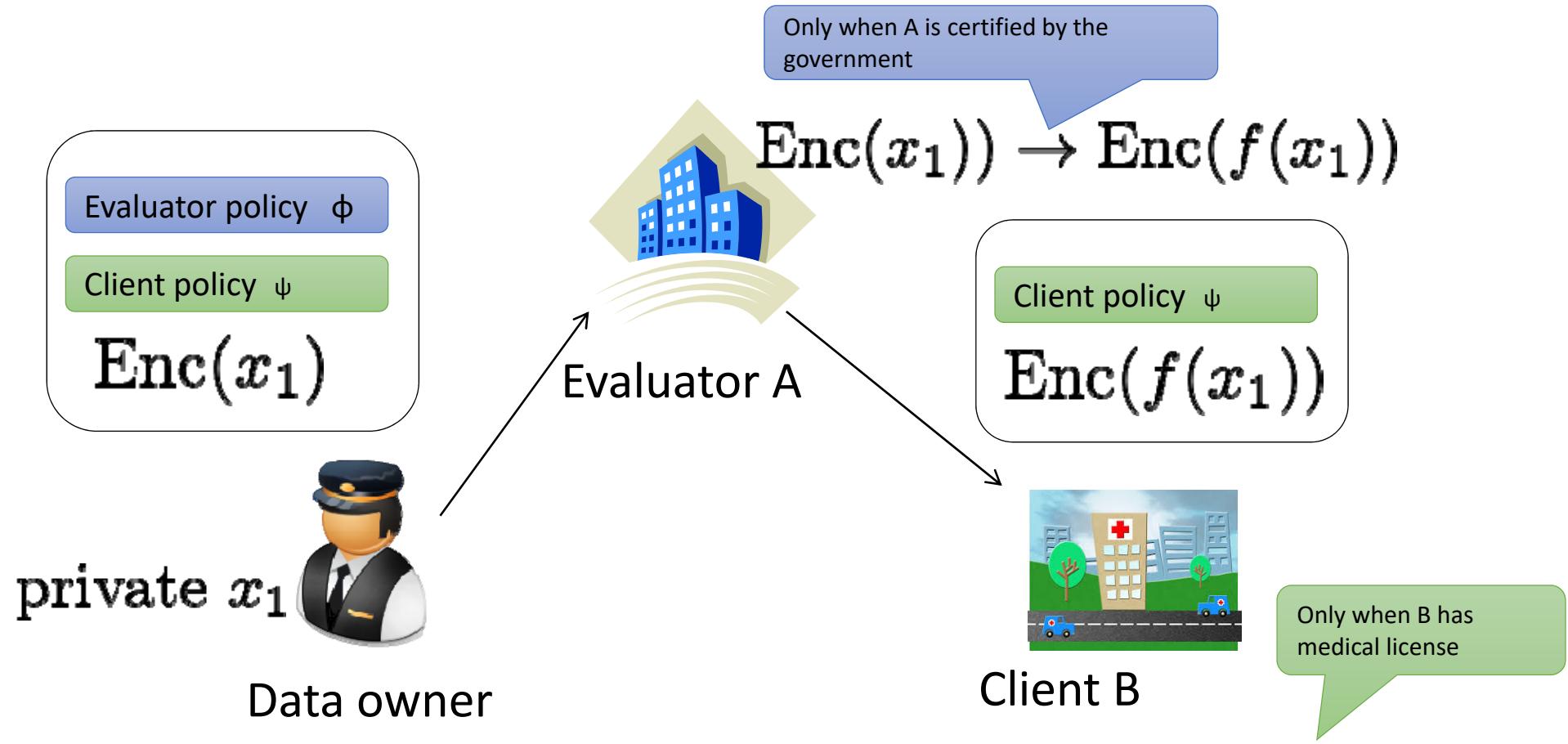
- ・個人ゲノムは解析は高コストだが生涯変化しない
- ・準同型暗号化個人ゲノムをクラウド上に蓄積
- ・様々な機関から様々な用途で安全に利用できるように個人ゲノム解析のサービス提供

# Secure Function Evaluation with Privacy Policy Enforcement

Sakuma, Lu, Nishide, Kunihiro, under review.



# Secure function evaluation with privacy policy enforcement



## Privacy policy

- Evaluator policy defines entities who are allowed to  $f(x)$
- Client policy defines entities who are allowed to obtain  $f(x)$

# Intuition of policy enforcement

Symmetric key enc. (SKE)



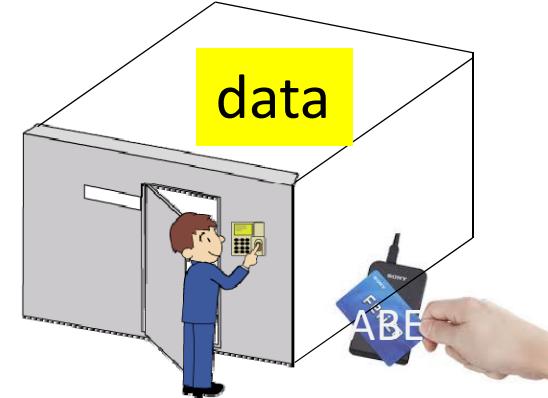
Those who have the key can see the data

Homomorphic enc. (HE)



Those who do not have the key cannot see the data, but can analyze the data

Attribute-based enc. (ABE)



Those who have a specified attribute can see the data

Privacy policy enforcement

Data encryption



Those who have the key can analyze the data

Encapsulated key



Those who have a specified attribute can obtain the key

Privacy policy enforcement scheme  
HE: $[x]$ , SKE: $\{x\}_K$ , ABE: $\langle x \rangle_\Psi$

CA KeyGeneration

**Evaluator  $\psi$  (not eligible)**

Can do nothing because  
 $\psi$  does not satisfy  $\Psi_E$

**Evaluator (eligible)**

DecABE:  $\langle K^E || K^C \rangle_\Psi \rightarrow K^E || K^C$

DecHE:  $\{[x_i]\}_K \rightarrow [x_i]$

EvalHE:  $[x_i] \rightarrow [f(x_i)]$

EnforceClientPol:  $[f(x_i)] \rightarrow [f(x_i)]_\Psi$

**Cloud**

Encrypted Data  
 $\{[x_i]\}_K^E$

Encapsulated Key  
 $\langle K^E || K^C \rangle_\Psi^E$

Encapsulated Result  
 $[f(x_i)]_\Psi^C$

DataEnc

$x_i \rightarrow [x_i] \rightarrow \{[x_i]\}_K^E$

KeyEncap

$: K^E || K^C \rightarrow \langle K^E || K^C \rangle_\Psi^E$

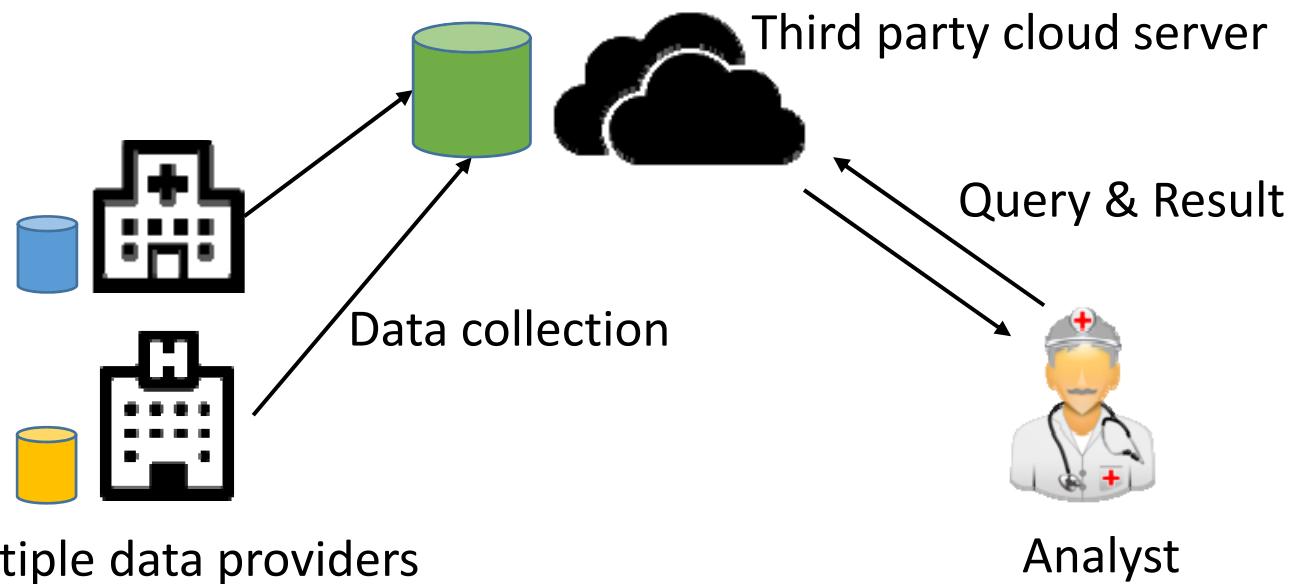
**Data contributor  $i$**

# Experimental results

	Disease risk prediction	SVM classification	Chi-square test																																																												
w/policy enforcement	<table border="1"> <caption>Evaluation cost of disease risk prediction</caption> <thead> <tr> <th>Dimension</th> <th>Without Enforcement (s)</th> <th>Enforcement (s)</th> </tr> </thead> <tbody> <tr><td>100</td><td>~2</td><td>~2</td></tr> <tr><td>200</td><td>~7</td><td>~7</td></tr> <tr><td>400</td><td>~11</td><td>~11</td></tr> <tr><td>800</td><td>~11</td><td>~11</td></tr> <tr><td>1600</td><td>~21</td><td>~21</td></tr> <tr><td>3200</td><td>~48</td><td>~48</td></tr> </tbody> </table>	Dimension	Without Enforcement (s)	Enforcement (s)	100	~2	~2	200	~7	~7	400	~11	~11	800	~11	~11	1600	~21	~21	3200	~48	~48	<table border="1"> <caption>Evaluation cost of SVM binary classification</caption> <thead> <tr> <th>#support vector</th> <th>Without Enforcement (s)</th> <th>Enforcement (s)</th> </tr> </thead> <tbody> <tr><td>5</td><td>~5</td><td>~5</td></tr> <tr><td>10</td><td>~10</td><td>~10</td></tr> <tr><td>15</td><td>~15</td><td>~15</td></tr> <tr><td>20</td><td>~20</td><td>~20</td></tr> <tr><td>25</td><td>~25</td><td>~25</td></tr> <tr><td>30</td><td>~30</td><td>~30</td></tr> </tbody> </table>	#support vector	Without Enforcement (s)	Enforcement (s)	5	~5	~5	10	~10	~10	15	~15	~15	20	~20	~20	25	~25	~25	30	~30	~30	<table border="1"> <caption>Evaluation cost of hypothesis testing</caption> <thead> <tr> <th>#subjects</th> <th>Without Enforcement (s)</th> <th>Enforcement (s)</th> </tr> </thead> <tbody> <tr><td>100</td><td>~30</td><td>~30</td></tr> <tr><td>200</td><td>~70</td><td>~70</td></tr> <tr><td>400</td><td>~120</td><td>~120</td></tr> <tr><td>800</td><td>~230</td><td>~230</td></tr> <tr><td>1600</td><td>~380</td><td>~100</td></tr> </tbody> </table>	#subjects	Without Enforcement (s)	Enforcement (s)	100	~30	~30	200	~70	~70	400	~120	~120	800	~230	~230	1600	~380	~100
Dimension	Without Enforcement (s)	Enforcement (s)																																																													
100	~2	~2																																																													
200	~7	~7																																																													
400	~11	~11																																																													
800	~11	~11																																																													
1600	~21	~21																																																													
3200	~48	~48																																																													
#support vector	Without Enforcement (s)	Enforcement (s)																																																													
5	~5	~5																																																													
10	~10	~10																																																													
15	~15	~15																																																													
20	~20	~20																																																													
25	~25	~25																																																													
30	~30	~30																																																													
#subjects	Without Enforcement (s)	Enforcement (s)																																																													
100	~30	~30																																																													
200	~70	~70																																																													
400	~120	~120																																																													
800	~230	~230																																																													
1600	~380	~100																																																													
wo/policy enforcement																																																															
task	$\text{risk} = \mathbf{w}^T \mathbf{x}$	$f(\mathbf{x}) = \sum_{\mathbf{x}_i^{SV} \in SV} \alpha_i k(\mathbf{x}, \mathbf{x}_i^{SV})$	$\chi^2 = \sum_{i=1}^4 \frac{(o_i - e_i)^2}{e_i}$																																																												
HE	Paillier	Ring-LWE	Ring-LWE																																																												
overhead	7%<	1.4%<	20%<																																																												

Privacy policy enforcement is not that costly!

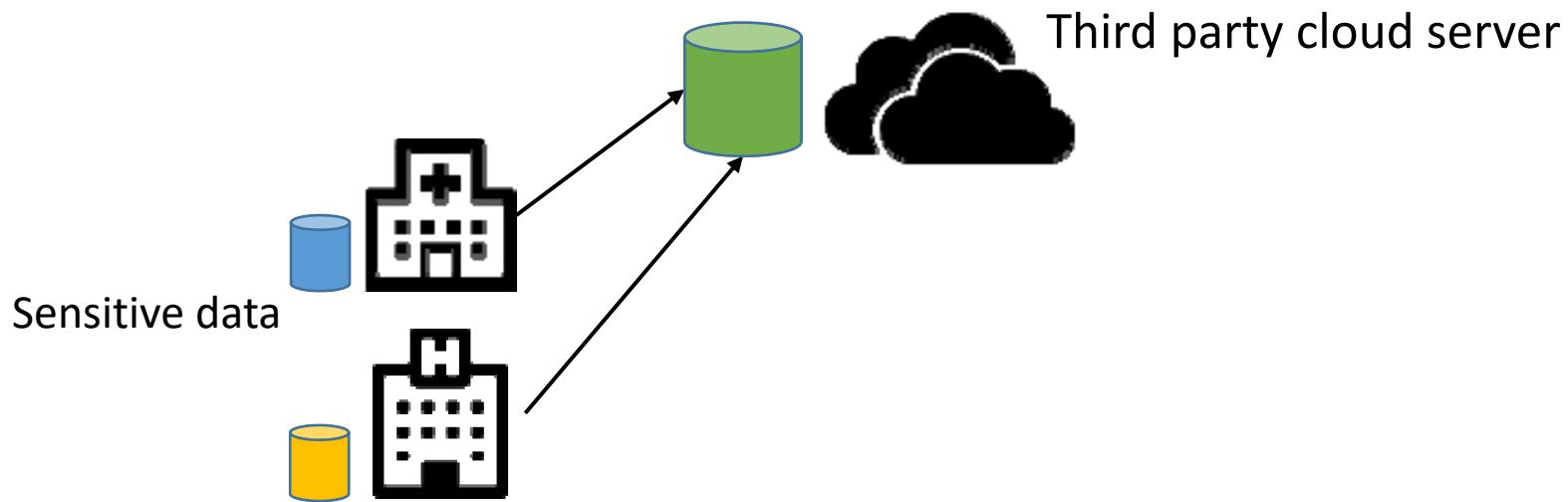
# Statistical Analysis on the Cloud



Cloud computing is useful for statistical analysis

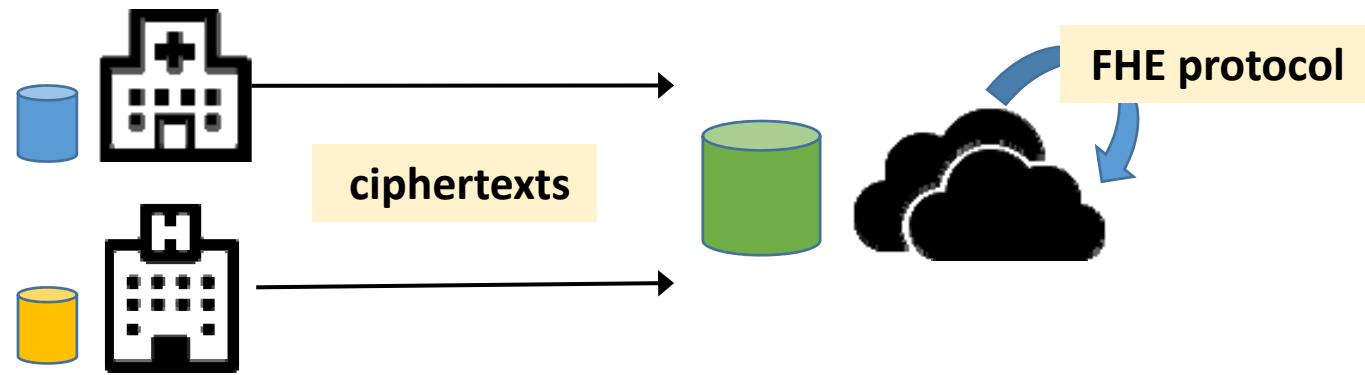
- Gather distributed data, and reduce hardware cost.
- Minimal interactions between data providers and the cloud.
- The cloud does most of the work for the analyst.

# Cloud Computing with Sensitive Data



- Using outside cloud servers raises privacy concerns.
  - E.g, medical records, federal data.
- We want to calculate statistics on the cloud while keeping the data secret.

# FHE on the Cloud Environment



- Less development cost
  - Suitable to outsourcing
- But might be inefficient in practice
  - Encrypt bits one by one.
  - 1~10 ms per evaluation.
  - 1~10 megabytes per ciphertext.

# Using Fully Homomorphic Encryption for Statistical Analysis of Categorical, Ordinal and Numerical Data

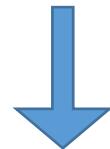
NDSS'17 Lu, Kawasaki, Sakuma

- Two new efficient FHE-based primitives:
  - *Matrix Operations*
  - *Batch Greater-than*
- Secure statistical protocols:
  - histogram (count),
  - order of counts,
  - contingency table (with cell-suppression),
  - percentile,
  - principal component analysis (PCA),
  - linear regression.
- Source codes: <https://github.com/fionser/CODA>

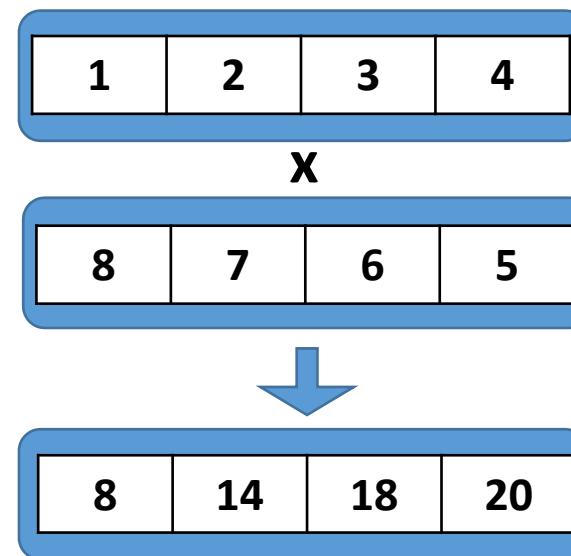
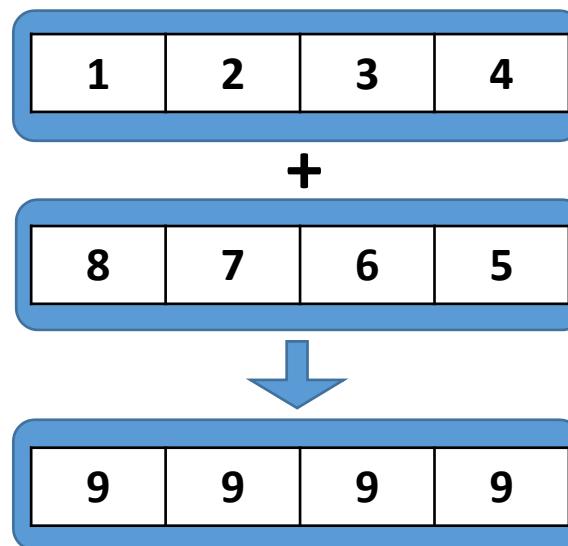
# Preliminaries: Packing (Batching)

- Enable to encrypt and process **vectors** at no extra cost.

Single  
homomorphic  
operation



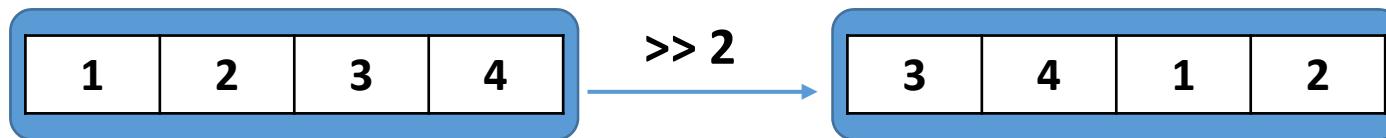
Multiple results



- Fewer ciphertexts
- Faster computation

# Preliminaries: Slot Manipulation

*Rotate* slots of the encrypted vector.



*Replicate* a specific slot.

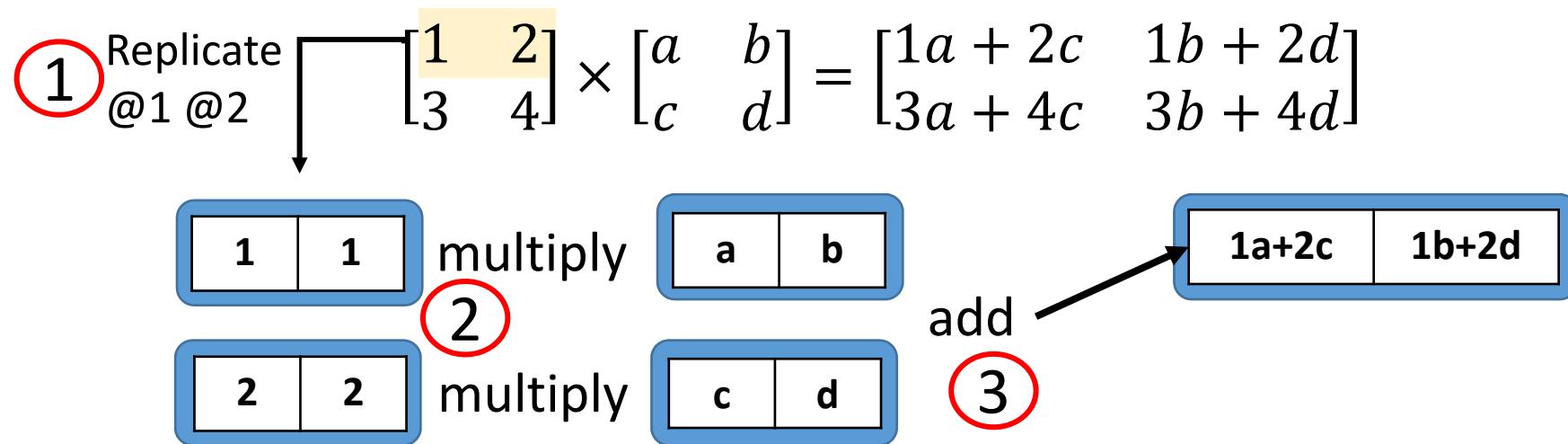


# Proposed Matrix Primitive

- Used for adding & multiplying encrypted matrices
- Encrypt each row separately by packing.
  - Row-wise encryption.
  - Horizontally partitioned data
- Efficient and layout consistent.
  - $O(N^2)$  homomorphic operations.

# Matrix Multiplication[1/2]

- Encrypt the matrix row by row with packing.

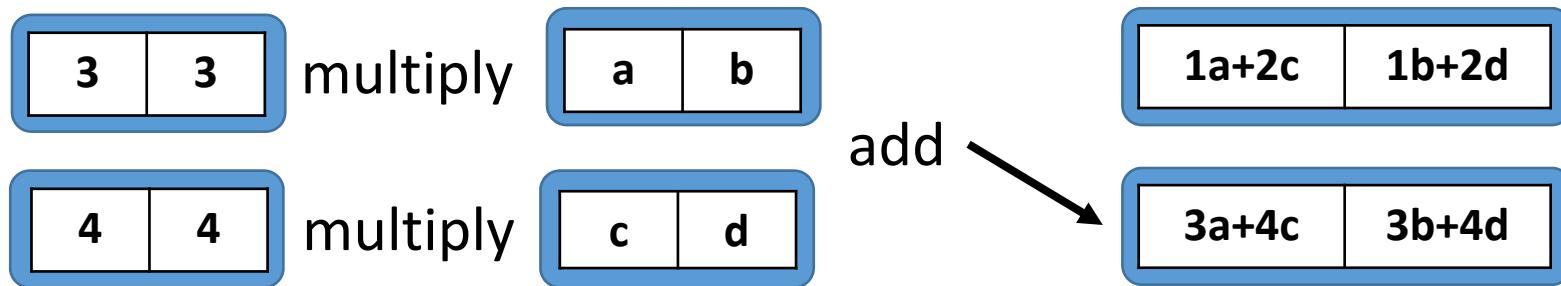


# Matrix Multiplication[1/2]

- Encrypt the matrix row by row with packing.

Replicate @1 @2

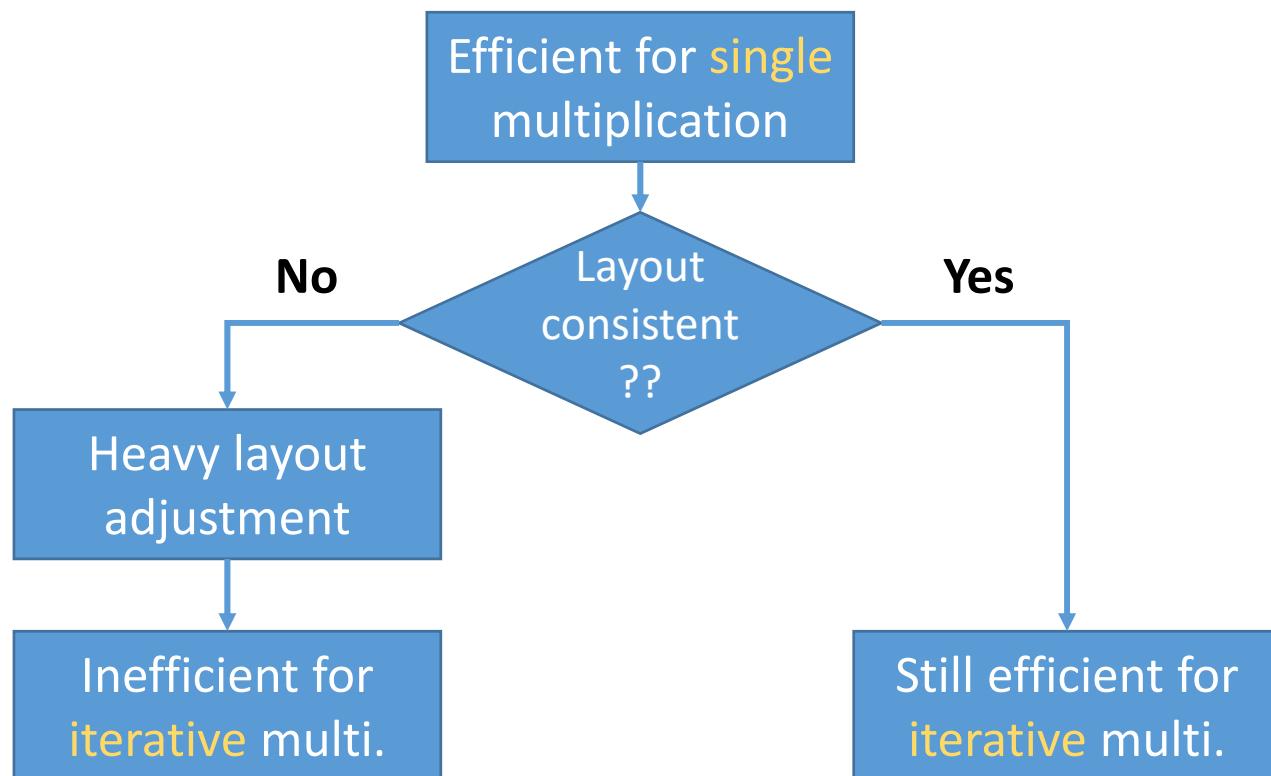
$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \times \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 1a + 2c & 1b + 2d \\ 3a + 4c & 3b + 4d \end{bmatrix}$$



- $N^2$  replications, multiplications and additions
  - $O(N^2)$  complexity compared to  $O(N^3)$  (no packing).
- Also row-wisely encrypted resulting matrix.

# Matrix Multiplication[2/2]

- Layout consistency is important for developing efficient statistical protocols.
  - Statistical algorithms need iterative matrix multiplications



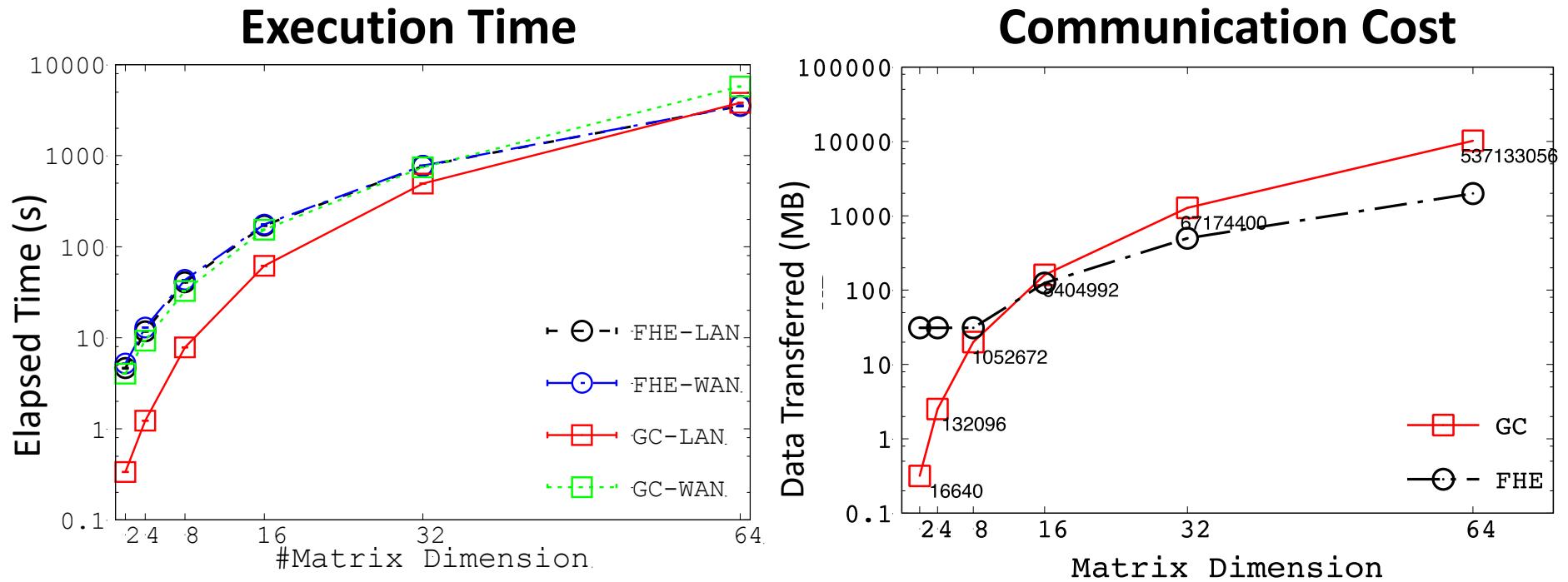
# Experimental Settings of Matrix Primitive

- Implementations:
  - FHE: HElib (C++ based)
  - GC : OblivM (java based)
- Evaluated on 32-bit integers
- Networks:
  - LAN (about 88 Mbps)
  - WAN (about 48 Mbps)

HElib. <https://github.com/shaih/HElib>.

Liu et al. *OblivM: A programming framework for secure computation*. 2015.

# Evaluation of Matrix Primitive



- When do iterative multiplications, FHE-based primitive can offer better performance.
  - Save communication cost between each iteration

# Greater-than (GT) Primitive

$$\text{GT}(e(x), e(y)) \rightarrow e(x >? y) \text{ s.t. } 0 \leq x, y \leq D$$

- [Golle06] based on Paillier cryptosystem:  
*if  $x > y$  then  $\exists k \in [1, D] \rightarrow x - y - k = 0$*
- Combination with packing gives great improvements:

$$e([x, \dots, x]) - e([y, \dots, y]) - [1, 2, \dots, D] \rightarrow e(\boldsymbol{\eta})$$

  
Replicated D times

- $0 \in \boldsymbol{\eta} \iff x > y$  (i.e., decryption is needed)
- Complexity from  $D$  to  $[D/\ell]$ .

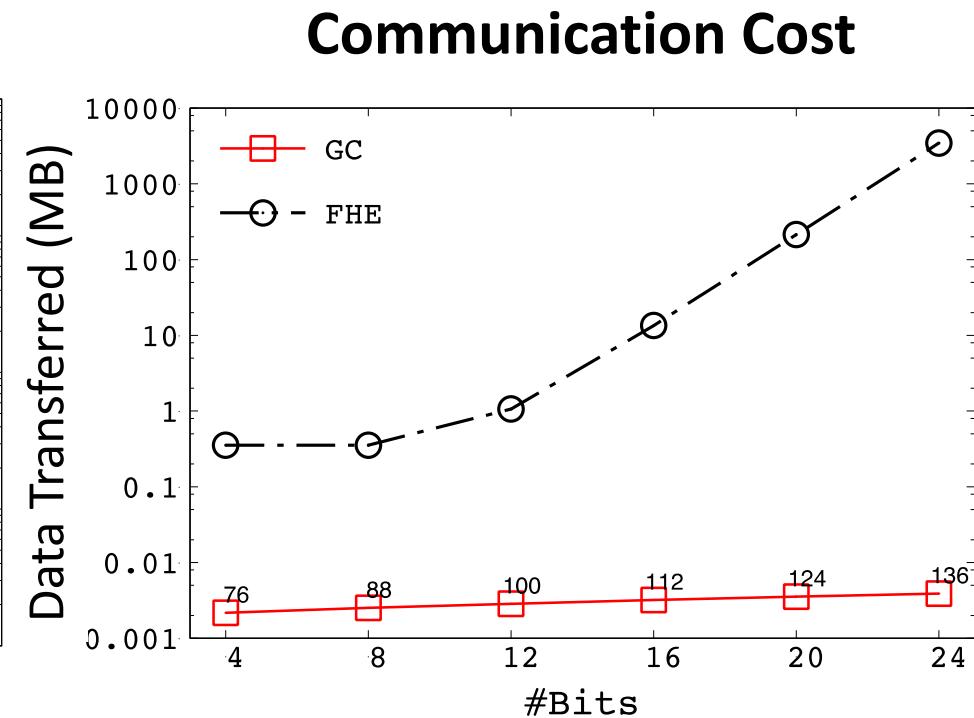
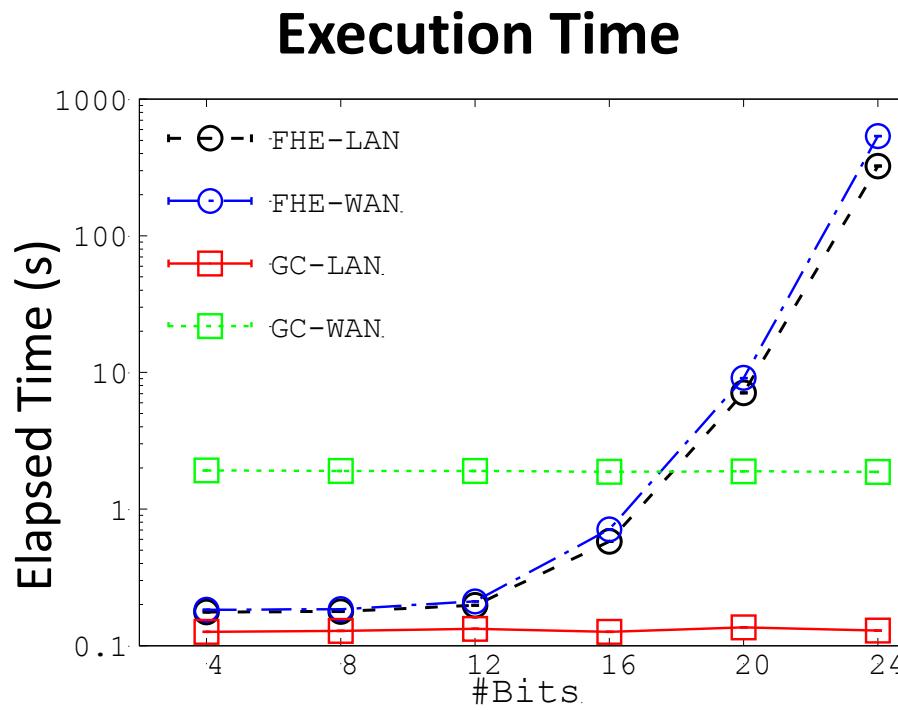
# Experimental Settings for GT Primitive

- Implementations:
  - FHE: HElib (C++ based)
  - GC : OblivM (java based)
- Domain  $D = 2^4 \sim 2^{24}$
- Number of slots  $\ell \approx 1700$ .
- Networks:
  - LAN (about 88 Mbps)
  - WAN (about 48 Mbps)

HElib. <https://github.com/shaih/HElib>.

Liu et al. *OblivM: A programming framework for secure computation*. 2015.

# Evaluation of Greater-than Primitive



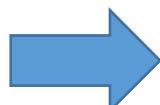
Works for small domains, which is enough for ordinal statistics.

# Secure Statistical Protocols

- Contingency table with cell-suppression protocol:
  - Use the greater-than primitive.
  - One round protocol between cloud and analyst.
- Linear regression protocol:
  - Use the matrix primitive.
  - Two rounds protocol.
  - Use a plaintext precision expansion technique (discuss it latter).

# Contingency Table

Gender	Smoke
Male	Smoker
Female	Non-smoker
Male	Non-Smoker



$$K_1 = 2$$

	Smoker	Non-smoker
Male	1	1
Female	0	1

Categorical data

$$K_2 = 2$$

Contingency Table

- Indicator encoding:

$$\text{Male} \rightarrow [1, 0], \quad \text{Female} \rightarrow [0, 1]$$

$$\text{Smoker} \rightarrow [1, 0], \text{Non-smoker} \rightarrow [0, 1]$$

- Basic Idea: **multiply & rotate**

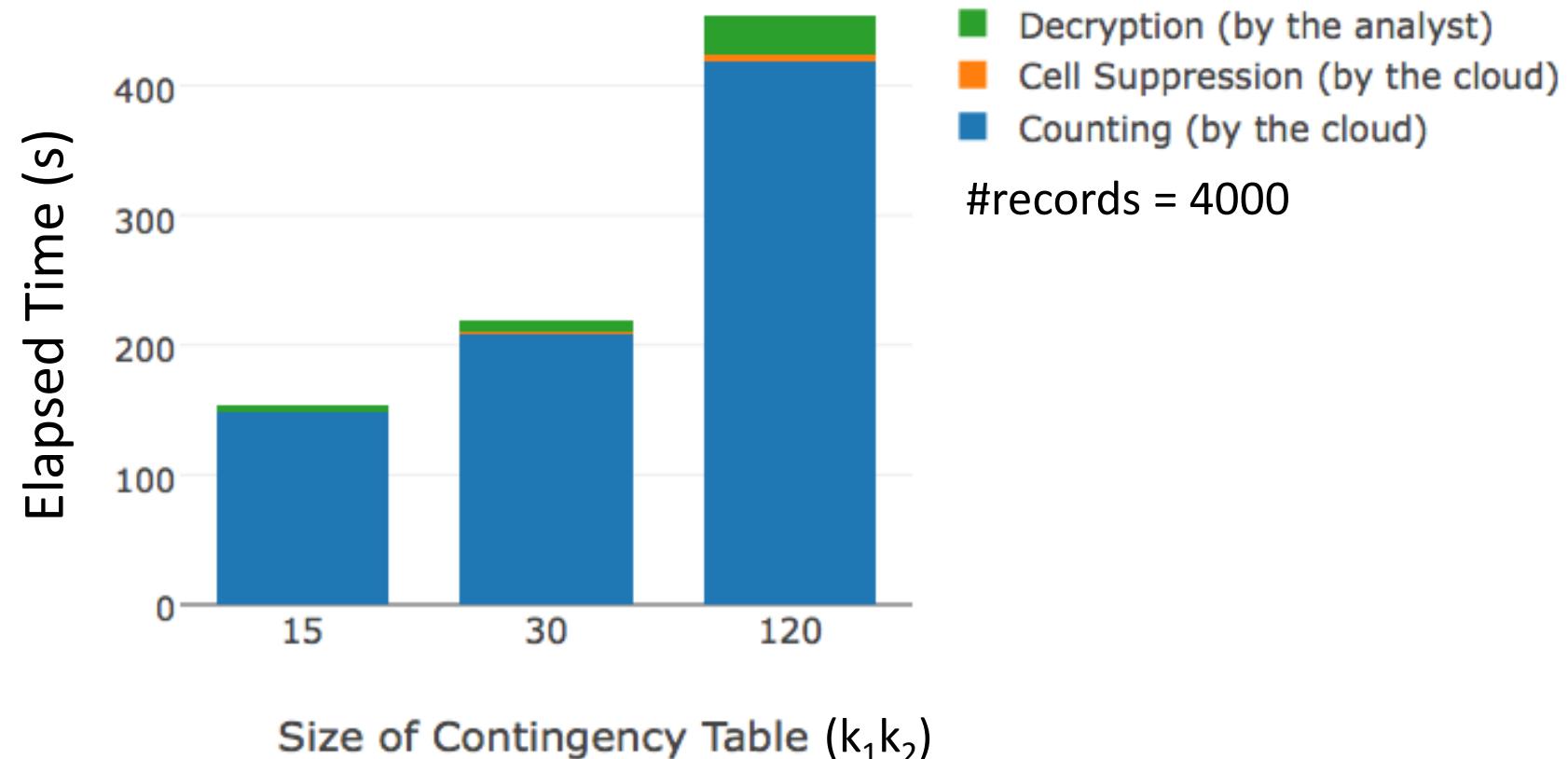
$[a_1, a_2] \times [b_1, b_2]$  counts Male-Smoker, and Female-Nonsmoker

$[a_1, a_2] \times ([b_1, b_2] \gg 1) = [a_1, a_2] \times [b_2, b_1]$  gives the other two counts.

- Improvement with no extra preprocessing

- $O(\max(k_1, k_2)) \Rightarrow O(\log k_1 k_2)$ .

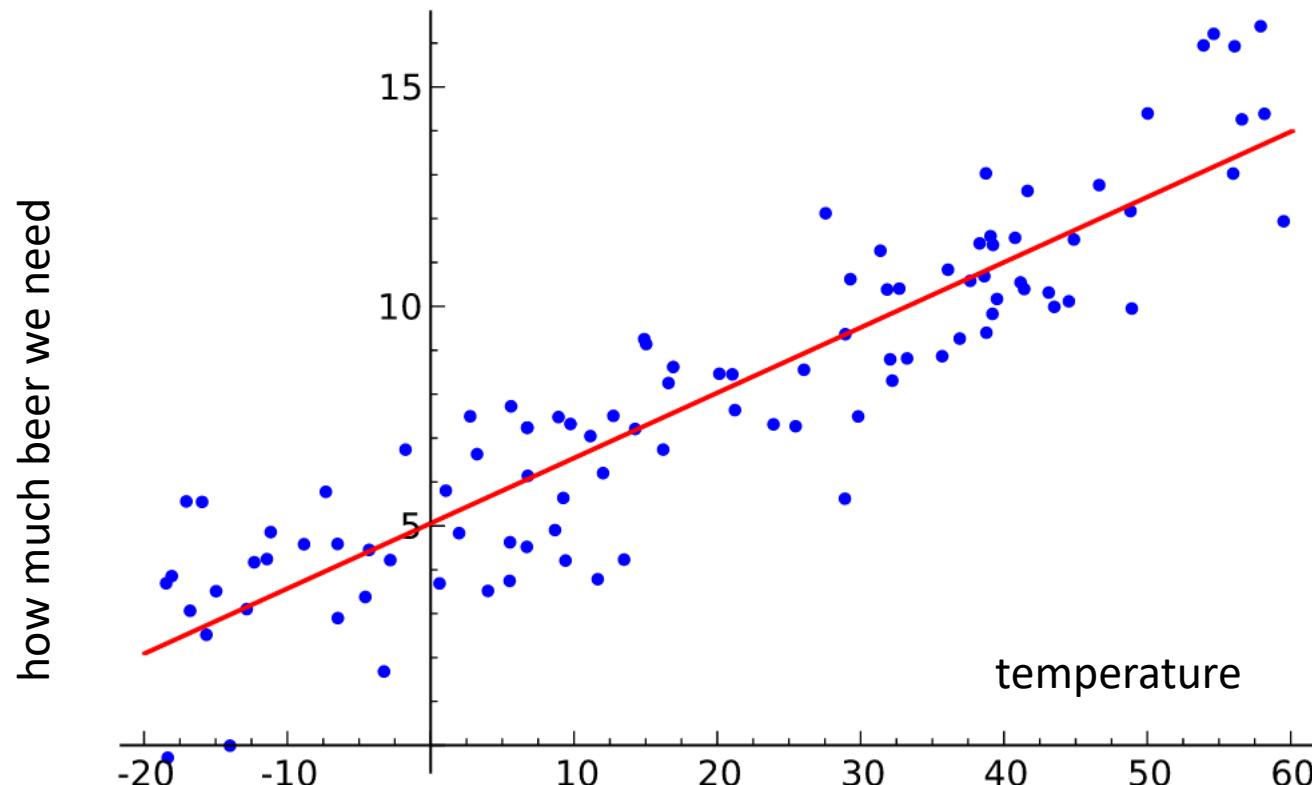
# Contingency Table Performance Evaluation



- Complexity increases logarithmically with the table sizes.
- Most of the work (>90%) done by the cloud.

# Liner regression

- From data  $\{(x_i, y_i)\}_i$ , computes a model  $w$  s.t.  
$$w = (X^T X)^{-1} X^T y$$



# How to solve LR over FHE

- The inversion of an encrypted matrix.

*Division-free Matrix Inversion ( $\mathbf{Q}, \lambda$ ):*

set  $\mathbf{A}^{(1)} = \mathbf{Q}$ ,  $\mathbf{R}^{(1)} = \mathbf{I}$ ,  $a^{(1)} = \lambda$ , and *iterate*

Layout consistency  
leads to efficient  
iterative protocols.

$$\mathbf{R}^{(t+1)} = 2a^{(t)}\mathbf{R}^{(t)} - \mathbf{R}^{(t)}\mathbf{A}^{(t)}$$

$$\mathbf{A}^{(t+1)} = 2a^{(t)}\mathbf{A}^{(t)} - \mathbf{A}^{(t)}\mathbf{A}^{(t)}$$

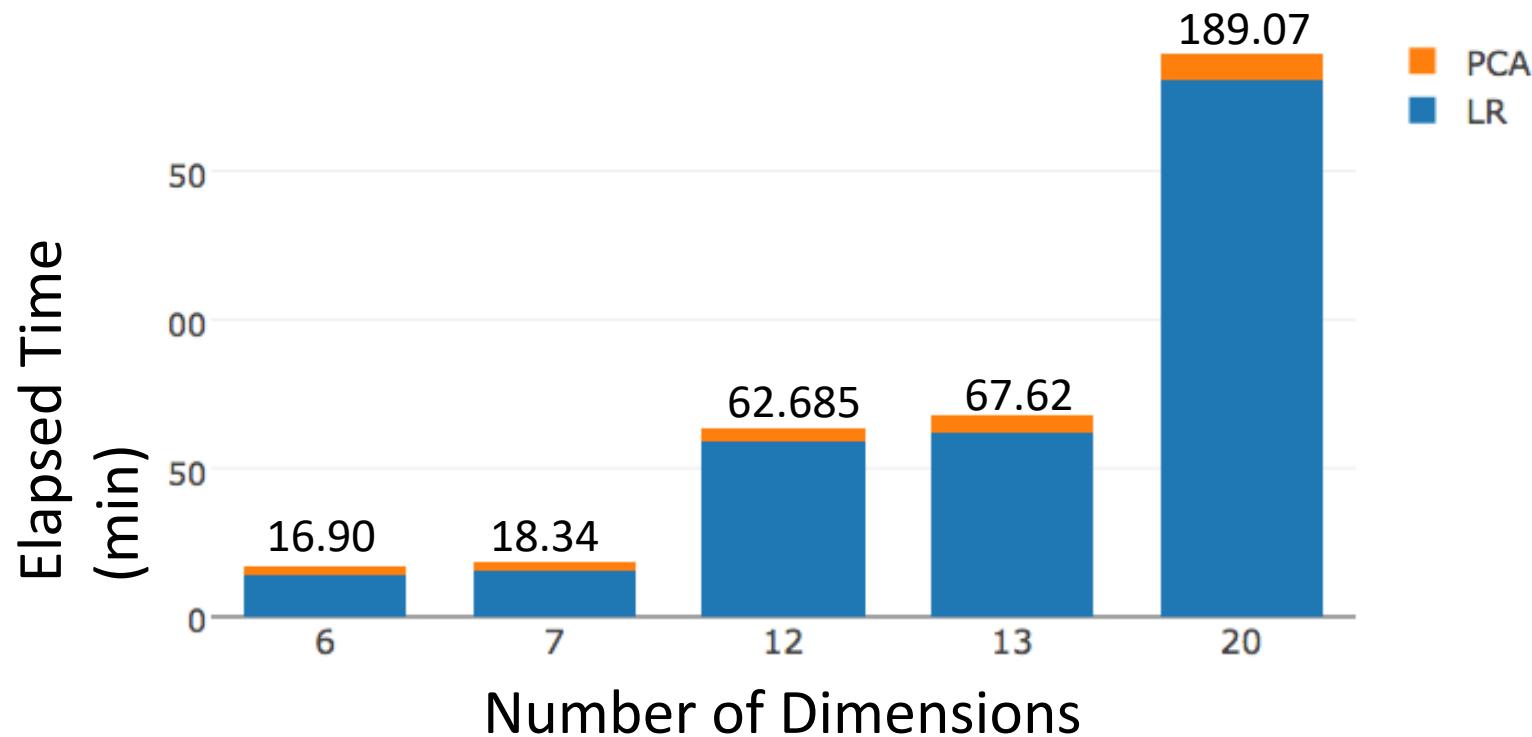
$$a^{(t+1)} = a^{(t)}a^{(t)}$$

[Guo06]  $\mathbf{R}^{(t)}$  gives a good *approximation* to  $\lambda^{2^t} \mathbf{Q}^{-1}$  if  $\lambda$  is close to largest eigenvalue of  $\mathbf{Q}$  (use PCA to compute  $\lambda$ ).

# Plaintext Precision Expansion (PPE)

- Division-free algorithms introduce large integers. ( $\lambda^{2^t}$ )
  - But the current FHE library allows at most 60-bit integers.
- Allows division-free algorithms without changing the FHE library.
- Uses  $K$  different FHE parameters (each  $b$ -bit  $< 60$ )
  - Achieves an equivalent  $\textcolor{red}{Kb}$ -bit parameter.
  - Increases the time by  $K$  times, but naturally parallelizable.
- Direct application of the Chinese Remainder Theorem.

# Experiments: Linear Regression



- Negligible decryption time (less than 2 s).
- 20x faster than previous FHE solution [Wu et al. 12]
  - 5 dimensions (400+ mins).
- Good scalability (reduced execution using more cores). 35

# Summary

- Secure statistical analysis in the cloud with multiple data providers.
- Two primitives
  - Matrix operation and greater-than
- Two protocols.
  - Contingency table and linear regression.
- Encoding and packing can improve FHE's balance between generality and efficiency.

# The Gap we see

- Epidemiology
  - Carefully design the objective of analysis
  - Collect a limited amount of well-controlled data
  - We evaluate whether or not the result is “significant”
- Big data analysis
  - Collect a huge amount of not well-controlled data
  - Objective of analysis is explorative
  - We can find a “good finding” if we are lucky
- How can we have “significant” finding with big data approaches?