# Homomorphic Encryption
# Based Secure Genome Data Analysis

Miran Kim[*] and Kristin Lauter[†]

[*]Seoul National University
[†]Microsoft Research

iDASH Privacy&Security Workshop, March 16, 2015

# Secure Outsourcing GWAS

# Minor Allele Frequency

- There are 200 people, and each of them has 311 genotypes.
- Each genotype has two kinds of SNPs.

- Data Encoding

  For a fixed genotype, suppose that 200 people have "AT, AT, AA, ..., TT". Then the encoding method is as follows:
  - If the pair consists of different SNPs (AT), then encode it into '1'.
  - The first pair with the same SNP (AA) is encoded into '0'.
  - Then the other one (TT) is encoded into '2'.

  ($\Rightarrow$ the encoded value means the number of 'T' in the individual SNPs.)

# Minor Allele Frequency

- There are 200 people, and each of them has 311 genotypes.
- Each genotype has two kinds of SNPs.

- Data Encoding

  For a fixed genotype, suppose that 200 people have "AT, AT, AA, ..., TT". Then the encoding method is as follows:
  - If the pair consists of different SNPs (AT), then encode it into '1'.
  - The first pair with the same SNP (AA) is encoded into '0'.
  - Then the other one (TT) is encoded into '2'.

  ($\Rightarrow$ the encoded value means the number of 'T' in the individual SNPs.)
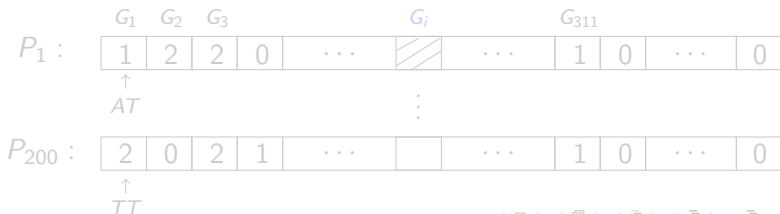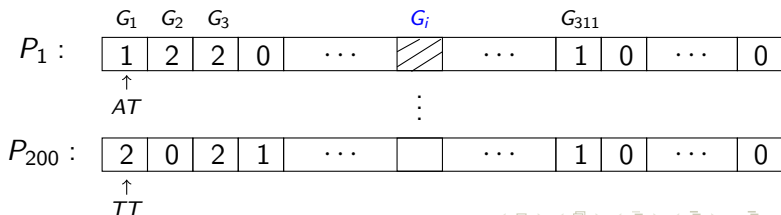
# Minor Allele Frequency

- There are 200 people, and each of them has 311 genotypes.
- Each genotype has two kinds of SNPs.

- Data Encoding

  For a fixed genotype, suppose that 200 people have "AT, AT, AA, ..., TT". Then the encoding method is as follows:
  - If the pair consists of different SNPs (AT), then encode it into '1'.
  - The first pair with the same SNP (AA) is encoded into '0'.
  - Then the other one (TT) is encoded into '2'.

  ($\Rightarrow$ the encoded value means the number of 'T' in the individual SNPs.)

# Minor Allele Frequency

- Encryption & Evaluation



(We can perform the aggregate operations simultaneously for all the genotypes.)

- Decryption
  - Decrypt the ciphertext "$\sum_{i=1}^{200} C_i$" with the secret key.
  - Let $\ell_i$ be the value in the $i$'th slot.

- Decoding
  - For $1 \le i \le 311$, if $\ell_i > 200$, then $\ell_i \leftarrow (400 - \ell_i)$.
  - The minor allele frequency of the genotype $G_i$ is $\left(\frac{\ell_i}{400}\right)$.

# Minor Allele Frequency

- Encryption & Evaluation



(We can perform the aggregate operations simultaneously for all the genotypes.)

- Decryption
  - Decrypt the ciphertext "$\sum_{i=1}^{200} C_i$" with the secret key.
  - Let $\ell_i$ be the value in the $i$'th slot.

- Decoding
  - For $1 \leq i \leq 311$, if $\ell_i > 200$, then $\ell_i \leftarrow (400 - \ell_i)$.
  - The minor allele frequency of the genotype $G_i$ is $\left(\frac{\ell_i}{400}\right)$.

# Chi-squared Test

- ## Data Encoding

    - For each genotype, encode the given SNPs of case group and control group.

* Note that the result of chi-squared test is

$$\frac{n(ad - bc)^2}{r \cdot s \cdot g \cdot k} = \frac{800\,(a(400 - c) - c(400 - a))^2}{400 \cdot 400 \cdot g \cdot k}$$

$$= \frac{800\,(a - c)^2}{(a + c)(800 - (a + c))}$$

where 'a' and 'c' are the allele counts of some SNP in case and control group.

# Chi-squared Test

- Data Encoding
  - For each genotype, encode the given SNPs of case group and control group.

\* Note that the result of chi-squared test is

$$
\frac{n(ad-bc)^2}{r \cdot s \cdot g \cdot k} = \frac{800\,(a(400-c)-c(400-a))^2}{400 \cdot 400 \cdot g \cdot k}
$$

$$
= \frac{800\,(a-c)^2}{(a+c)(800-(a+c))}
$$

where '$a$' and '$c$' are the allele counts of some SNP in case and control group.

# Chi-squared Test

- ## Evaluation

  Let us denote $C_i$ and $C_i'$ the ciphertexts for the case&control groups.

  - Evaluate $\sum_{i=1}^{200} C_i$ ( $\stackrel{let}{=} \mathfrak{C}_{case}$) and $\sum_{i=1}^{200} C_i'$ ( $\stackrel{let}{=} \mathfrak{C}_{cont}$).
  - Compute "$\mathfrak{C}_{case} - \mathfrak{C}_{cont}$" and "$\mathfrak{C}_{case} + \mathfrak{C}_{cont}$"

- ## Decryption

  For the message space $\mathbb{Z}_t = [0, t)$,

  - den $\stackrel{let}{=} \text{Dec}(\mathfrak{C}_{case} + \mathfrak{C}_{cont}) = a + c \ (< t)$
  - num $\stackrel{let}{=} \text{Dec}(\mathfrak{C}_{case} - \mathfrak{C}_{cont}) = \begin{cases} a - c & \text{if } a > c, \\ (a - c) + t & \text{otherwise.} \end{cases}$

- ## Decoding

  - If num $> \frac{t}{2}$, then num $\leftarrow$ (num $- t$).
  - The result of chi-squared test is $\frac{800(\text{num})^2}{(\text{den})(800 - \text{den})}$

# Chi-squared Test

- Evaluation

  Let us denote $C_i$ and $C_i'$ the ciphertexts for the case&control groups.

  - Evaluate $\sum_{i=1}^{200} C_i$ ( $\overset{let}{=} \mathfrak{C}_{case}$ ) and $\sum_{i=1}^{200} C_i'$ ( $\overset{let}{=} \mathfrak{C}_{cont}$ ).
  - Compute "$\mathfrak{C}_{case} - \mathfrak{C}_{cont}$" and "$\mathfrak{C}_{case} + \mathfrak{C}_{cont}$"

- Decryption

  For the message space $\mathbb{Z}_t = [0, t)$,

  - $\text{den} \overset{let}{=} \text{Dec}(\mathfrak{C}_{case} + \mathfrak{C}_{cont}) = a + c \ (< t)$
  - $\text{num} \overset{let}{=} \text{Dec}(\mathfrak{C}_{case} - \mathfrak{C}_{cont}) = \begin{cases} a - c & \text{if } a > c, \\ (a - c) + t & \text{otherwise.} \end{cases}$

- Decoding

  - If $\text{num} > \frac{t}{2}$, then $\text{num} \leftarrow (\text{num} - t)$.
  - The result of chi-squared test is $\frac{800(\text{num})^2}{(\text{den})(800 - \text{den})}$

Secure Comparison

between Genomic Data

# Hamming Distance

- Two individuals have genotypes over many SNPs. For a fixed genotype,

$$d = \begin{cases} 1 & \text{if } (S_1 = \text{null}) \mathbin{||} (S_2 = \text{null}) \mathbin{||} (S_1.\text{alt} \neq S_2.\text{alt}) \\ 0 & \text{otherwise} \end{cases}$$

- $x[j] \overset{let}{=} j$-th bit of $x$, starting with the least significant bit of $x$.
  $\oplus : \mathrm{XOR}$ gate $(= \text{Add over } \mathbb{Z}_2)$, $\wedge : \mathrm{AND}$ gate $(= \text{Mult over } \mathbb{Z}_2)$.

| SVTYPE | $d$ |
|---|---|
| $SV_1$ or $SV_2 = \text{INS/DEL}$ | 0 |
| $SV_1$ or $SV_2 = \text{null}$ | 1 |
| $SV_1$ and $SV_2 = \text{SNP/SUB}$ | $\mathrm{EQU}(S_1, S_2) \oplus 1$ |

where $\mathrm{EQU}(S_1, S_2) = \begin{cases} 1 & \text{if } S_1 = S_2 \\ 0 & \text{o.w,} \end{cases} = \wedge_{j=1}^{\mu} \left( S_1[j] \oplus S_2[j] \oplus 1 \right)$

- We need the encodings to determine 'null' and 'INS/DEL'.

# Hamming Distance

- Two individuals have genotypes over many SNPs. For a fixed genotype,

$$d = \begin{cases} 1 & \text{if } (S_1 = \text{null}) \;||\; (S_2 = \text{null}) \;||\; (S_1.\text{alt} \neq S_2.\text{alt}) \\ 0 & \text{otherwise} \end{cases}$$

- $x[j] \stackrel{let}{=} j$-th bit of $x$, starting with the least significant bit of $x$.
  $\oplus : \mathrm{XOR}$ gate (= Add over $\mathbb{Z}_2$), $\wedge : \mathrm{AND}$ gate (= Mult over $\mathbb{Z}_2$).

| SVTYPE | $d$ |
|---|---|
| $SV_1$ or $SV_2 = $ INS/DEL | 0 |
| $SV_1$ or $SV_2 = $ null | 1 |
| $SV_1$ and $SV_2 = $ SNP/SUB | $\texttt{EQU}(S_1, S_2) \oplus 1$ |

where $\texttt{EQU}(S_1, S_2) = \begin{cases} 1 & \text{if } S_1 = S_2 \\ 0 & \text{o.w,} \end{cases} = \wedge_{j=1}^{\mu} \left( S_1[j] \oplus S_2[j] \oplus 1 \right)$

- We need the encodings to determine 'null' and 'INS/DEL'.

# Hamming Distance

- Data Encoding
  - Clean two datasets using POS, then make the merged list $L$.

  - For $i \in [1, \#(L)]$,

    define $m_i = \begin{cases} 1 & \text{if } POS_i \in L \\ 0 & \text{otherwise} \end{cases}$ and $h_i = \begin{cases} 0 & \text{if } SV_i = INS/DEL \\ 1 & \text{otherwise} \end{cases}$

    $\Rightarrow (m_1 \oplus m_2) = 1$ iff $(SV_{1,i} = \text{null})$ or $(SV_{2,i} = \text{null})$

    $(h_1 \wedge h_2) = 0$ iff $(SV_{1,i} = INS/DEL)$ or $(SV_{2,i} = INS/DEL)$

  - Encode the SNP string as follows:

    $$A \rightarrow 00, \ G \rightarrow 01, \ C \rightarrow 10, \ T \rightarrow 11,$$

    - Each SNP is encoded and concatenated each other.
    - Pad '1' at the end of the string, and '0' to make 21 bit string, say $S_i$.
    - In the case of missing genotype, it is encoded as '0' string.
    - For example, '$GTA$' is encoded as '01||11||00||1$\underbrace{0 \ldots 00}_{14}$'.

# Hamming Distance

- ### Data Encoding

  - Clean two datasets using POS, then make the merged list $L$.

  - For $i \in [1, \#(L)]$,

    define $m_i = \begin{cases} 1 & \text{if POS}_i \in L \\ 0 & \text{otherwise} \end{cases}$ and $h_i = \begin{cases} 0 & \text{if SV}_i = \text{INS/DEL} \\ 1 & \text{otherwise} \end{cases}$

    $\Rightarrow (m_1 \oplus m_2) = 1$ iff $(\text{SV}_{1,i} = \text{null})$ or $(\text{SV}_{2,i} = \text{null})$

    $\quad (h_1 \wedge h_2) = 0 \quad$ iff $(\text{SV}_{1,i} = \text{INS/DEL})$ or $(\text{SV}_{2,i} = \text{INS/DEL})$
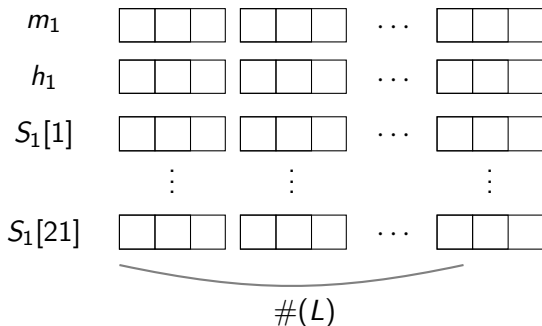
  - Encode the SNP string as follows:

    $$A \rightarrow 00, G \rightarrow 01, C \rightarrow 10, T \rightarrow 11,$$

    - ⋆ Each SNP is encoded and concatenated each other.
    - ⋆ Pad '1' at the end of the string, and '0' to make 21 bit string, say $S_i$.
    - ⋆ In the case of missing genotype, it is encoded as '0' string.
    - ⋆ For example, '$GTA$' is encoded as '$01||11||00||1\underbrace{0\ldots00}_{14}$'.

# Hamming Distance

- Encryption
  - Embed the data of $P_1(= m_{1,i}, h_{1,i}, S_{1,i})$ and $P_2$ $(= m_{2,i}, h_{2,i}, S_{2,i})$ into the plaintext slots in a bit-by-bit manner.
  - Encrypt the slots with the public key.

# Hamming Distance

- Evaluation
  - Evaluate the following binary circuit over encrypted data:
  
  $$(h_{1,i} \wedge h_{2,i}) \wedge \left( (m_{1,i} \oplus m_{2,i}) \oplus \left( m_{1,i} \oplus m_{2,i} \oplus 1 \right) \wedge \left( \text{EQU}(S_{1,i}, S_{2,i}) \oplus 1 \right) \right)$$
  
  - Take $m = 8191$ so that we can embed 630 messages into one ciphertext and perform the operations simultaneously for all the messages.

- Decryption
  - Decrypt the evaluated value and let $\ell_i$ the value in the $i$'th slot.

- Decoding
  - Note that $\ell_i$ is the Hamming distance result of $i$'th genotype.
  - Compute $\sum_{i=1}^{\#(L)} \ell_i$.

# Edit Distance

- For each genotype, we let

$$n = \begin{cases} \text{len}(S.\text{alt}) & \text{if SV} = \text{SNP/SUB/INS} \\ \text{len}(S.\text{ref}) & \text{if SV} = \text{DEL} \end{cases}$$

- $d = \begin{cases} 0 & \text{if } (S_1.\text{ref} = S_2.\text{ref}) \text{ \& } (S_1.\text{alt} = S_2.\text{alt}) \\ \max(n_1, n_2) & \text{otherwise} \end{cases}$

| SVTYPE | $d$ |
|---|---|
| $(\text{SV}_1 = \text{INS}, \text{SV}_2 \neq \text{INS}) \| (\text{SV}_1 \neq \text{INS}, \text{SV}_2 = \text{INS})$ | $\max(n_1, n_2)$ |
| Otherwise | $\max(n_1, n_2) \wedge (\text{EQU}(S_1.\text{alt}, S_2.\text{alt}) \oplus 1)$ |

- ▸ We don't need the reference comparison anymore.
- ▸ We need an encoding which determine whether the genotype is INS or not.

# Edit Distance

- For each genotype, we let

$$n = \begin{cases} \text{len}(S.\text{alt}) & \text{if } SV = SNP/SUB/INS \\ \text{len}(S.\text{ref}) & \text{if } SV = DEL \end{cases}$$

- $d = \begin{cases} 0 & \text{if } (S_1.\text{ref} = S_2.\text{ref}) \ \& \ (S_1.\text{alt} = S_2.\text{alt}) \\ \max(n_1, n_2) & \text{otherwise} \end{cases}$

| SVTYPE | $d$ |
|---|---|
| $(SV_1 = INS, SV_2 \neq INS) \| (SV_1 \neq INS, SV_2 = INS)$ | $\max(n_1, n_2)$ |
| Otherwise | $\max(n_1, n_2) \wedge (\text{EQU}(S_1.\text{alt}, S_2.\text{alt}) \oplus 1)$ |

- ► We don't need the reference comparison anymore.
- ► We need an encoding which determine whether the genotype is INS or not.

# Edit Distance

- Data Encoding
  - Clean two datasets using POS, then make the merged list $L$.
  - For $i \in [1, \#(L)]$, define $e_i = \begin{cases} 1 & \text{if } SV_i = INS, \\ 0 & \text{o.w.} \end{cases}$

    $\Rightarrow (e_1 \oplus e_2 \oplus 1) = 1$ iff $((SV_{1,i} = I, SV_{2,i} \neq I)$ or $(SV_{1,i} \neq I, SV_{2,i} = I))$
  - Encode the SNP string as $S_i$. (The missing genotype is encoded as '0')
  - Encode the length of SNP string, say $n_i$.

- Encryption
  - Embed the data of $P_1 (= e_{1,i}, S_{1,i}, n_{1,i})$ and $P_2 (= e_{2,i}, S_{2,i}, n_{2,i})$ into the plaintext slots in a bit-by-bit manner.
  - Encrypt the slots with the public key.

# Edit Distance

- Evaluation

  For $\mu$-bit integer x and y,

  - $\texttt{C}(x, y) = \begin{cases} 1 & \text{if } x < y \\ 0 & \text{o.w.,} \end{cases} = c_\mu$

    - $c_1 = (1 \oplus x[1]) \wedge y[1]$,

    - $c_j = \big((1 \oplus x[j]) \wedge y[j]\big) \oplus \big((1 \oplus x[j] \oplus y[j]) \wedge c_{j-1}\big)$ for $2 \le j \le \mu$

  - $\texttt{max}(x, y)[j] = \begin{cases} y[j] & \text{if } x < y, \\ x[j] & \text{o.w.,} \end{cases}$

    $= \big((1 \oplus \texttt{C}(x, y)) \wedge x[j]\big) \oplus \big(\texttt{C}(x, y) \wedge y[j]\big)$
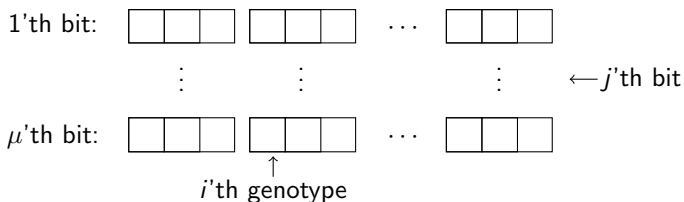
  - For $i \in [1, \#(L)]$ and $j \in [1, \mu]$, evaluate the circuits homomorphically:

    $$\Big(\big(\texttt{EQU}(S_{1,i}, S_{2,i}) \wedge (e_{1,i} \oplus e_{2,i} \oplus 1)\big) \oplus 1\Big) \wedge \texttt{max}(n_{1,i}, n_{2,i})[j]$$

# Edit Distance

- **Decryption**
  - ▸ Decrypt the evaluated values, and let $\ell_{i,j}$ the $i$'th value in the $j$'th slot.



$i$'th genotype

- **Decoding**
  - ▸ $\ell_i \overset{let}{=} \sum_{j=1}^{\mu} \ell_{i,j} \cdot 2^{j-1}$ is the Edit distance result of $i$'th genotype.
  - ▸ Compute $\sum_{i=1}^{\#(L)} \ell_i$.