

Low Latency Privacy Preserving Inference

Alon Brutzkuz, Oren Elisha, Ran Gilad-Bachrach

Microsoft Research, June 2019

Proceedings of the 36th International Conference on Machine Learning,
PMLR

Outline

- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

Outline

- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

Homomorphic encryption

- $D(E(m_1) * E(m_2)) = m_1 * m_2$
- $D(E(m_1) + E(m_2)) = m_1 + m_2$
- Partially homomorphic encryption
 - Either addition or multiplication
- Somewhat homomorphic encryption
 - Addition and multiplication with limited number of operations
- Fully homomorphic encryption
 - Addition and multiplication without limitation

Example

- $D(E(m_1) * E(m_2)) = m_1 * m_2$
- $D(E(m_1) + E(m_2)) = m_1 + m_2$



Example

- $D(E(m_1) * E(m_2)) = m_1 * m_2$
- $D(E(m_1) + E(m_2)) = m_1 + m_2$



Somewhat homomorphic encryption

- $D(E(m_1) * E(m_2)) = m_1 * m_2$
- $D(E(m_1) + E(m_2)) = m_1 + m_2$
- With limited number of operations



Somewhat homomorphic encryption

- $D(E(m_1) * E(m_2)) = m_1 * m_2$
- $D(E(m_1) + E(m_2)) = m_1 + m_2$

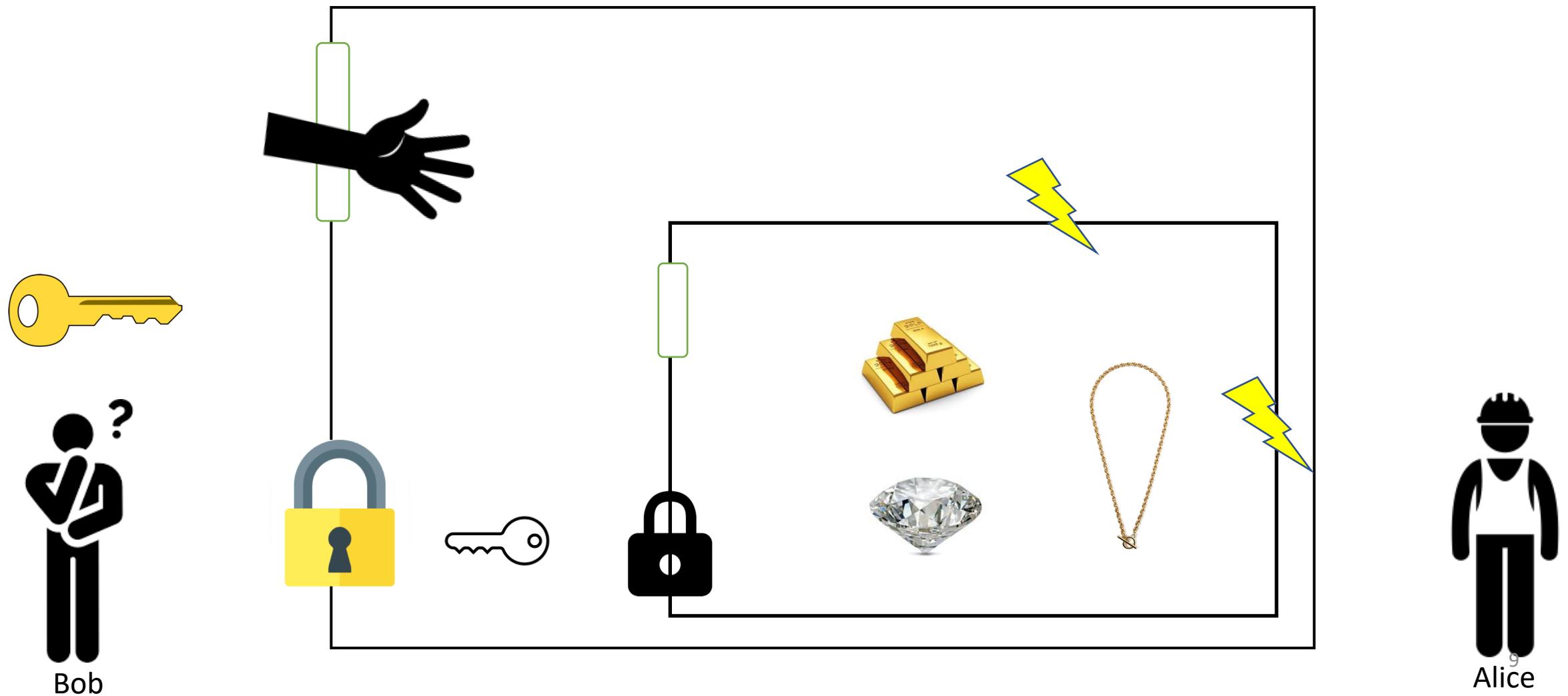


Addition

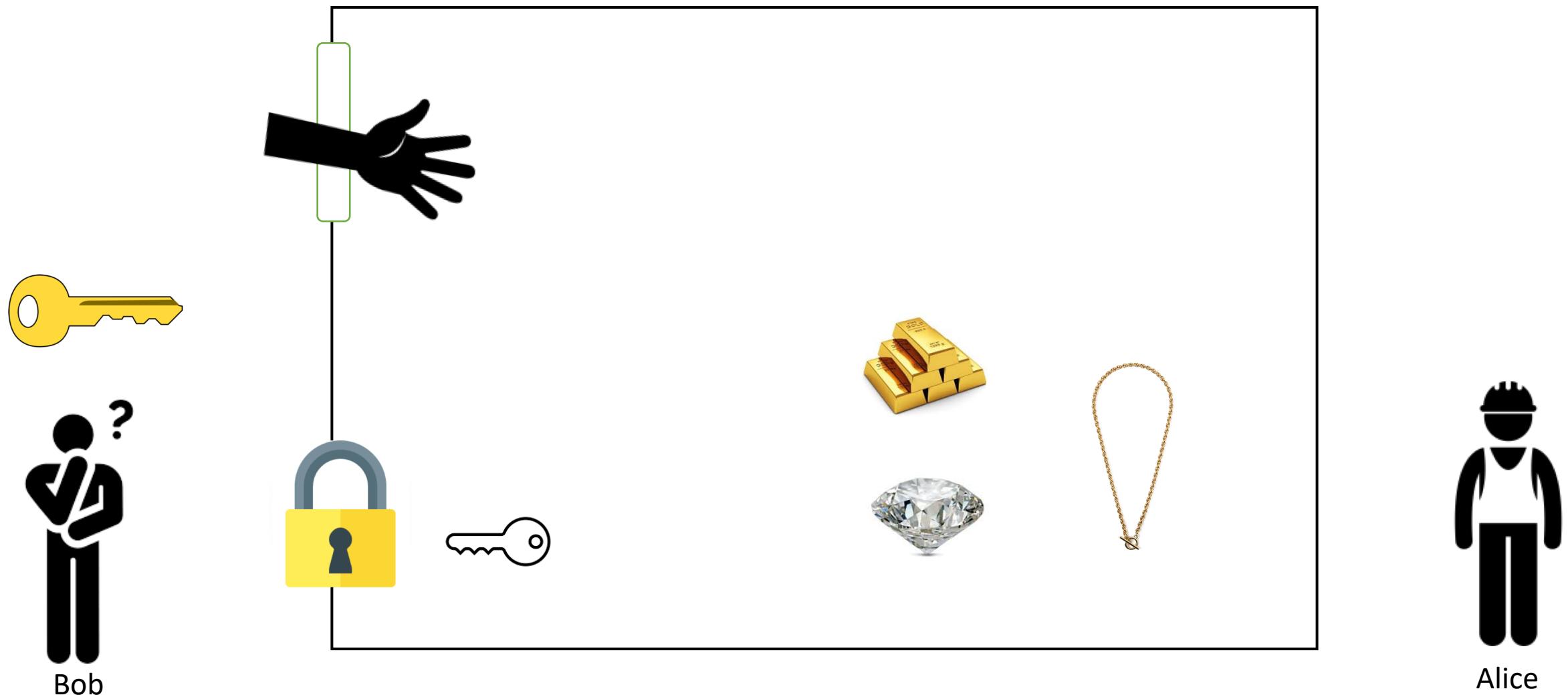


Multiplication

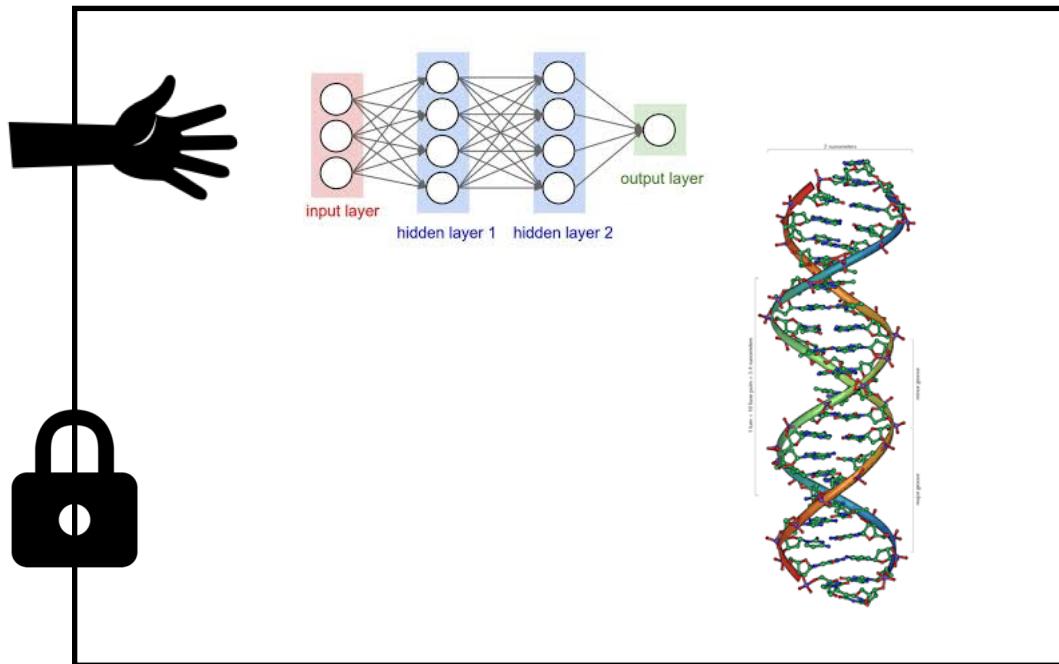
Bootstrapping



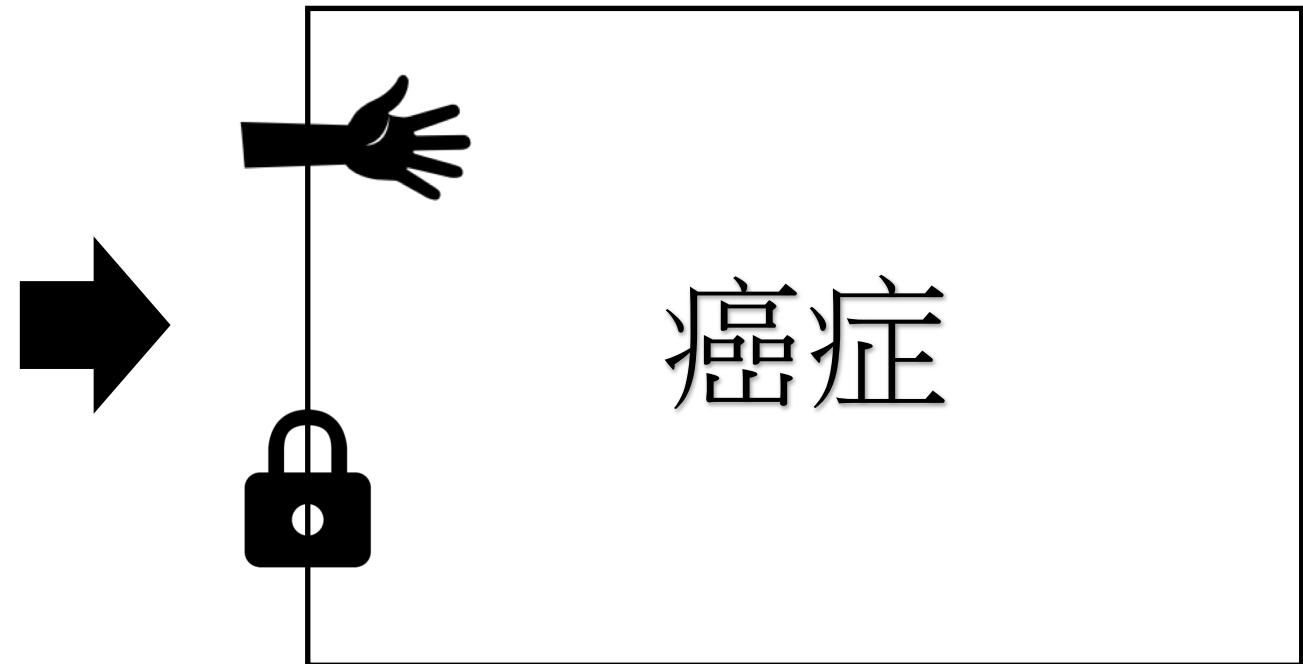
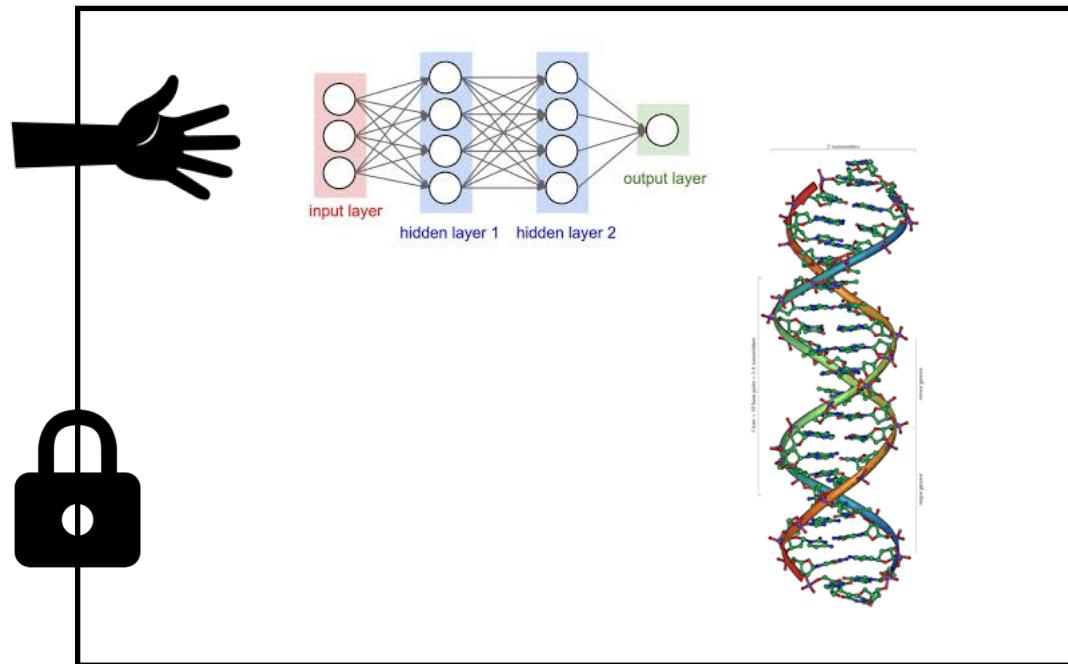
Fully homomorphic encryption



Analogy to encrypted inference



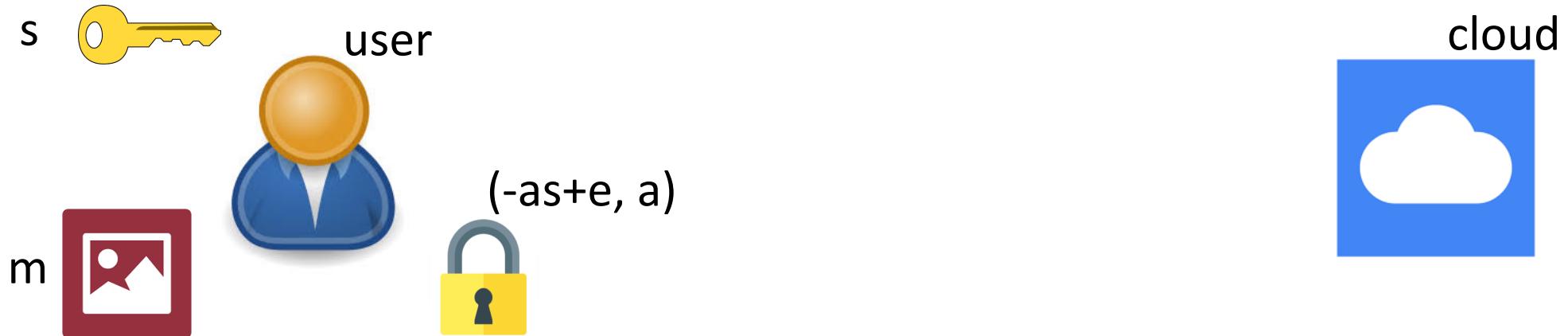
Analogy to encrypted inference



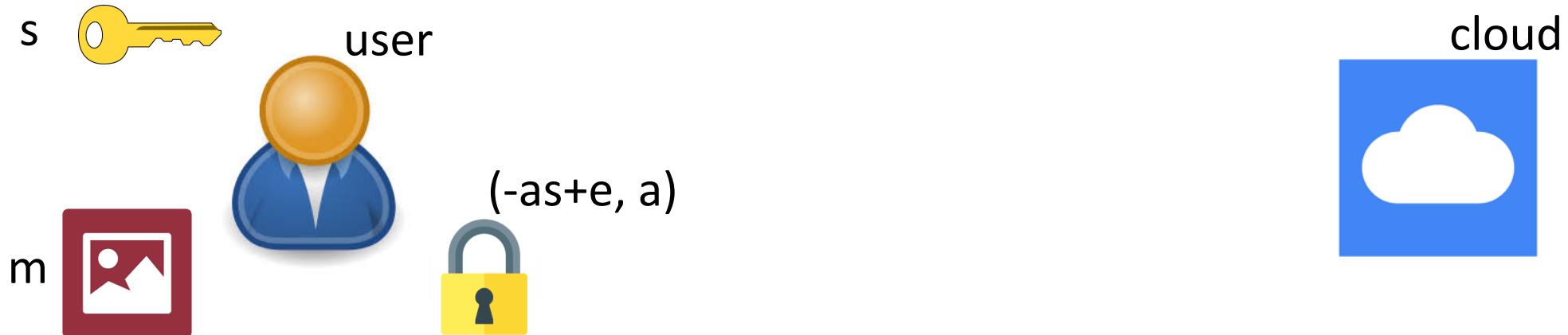
Outline

- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

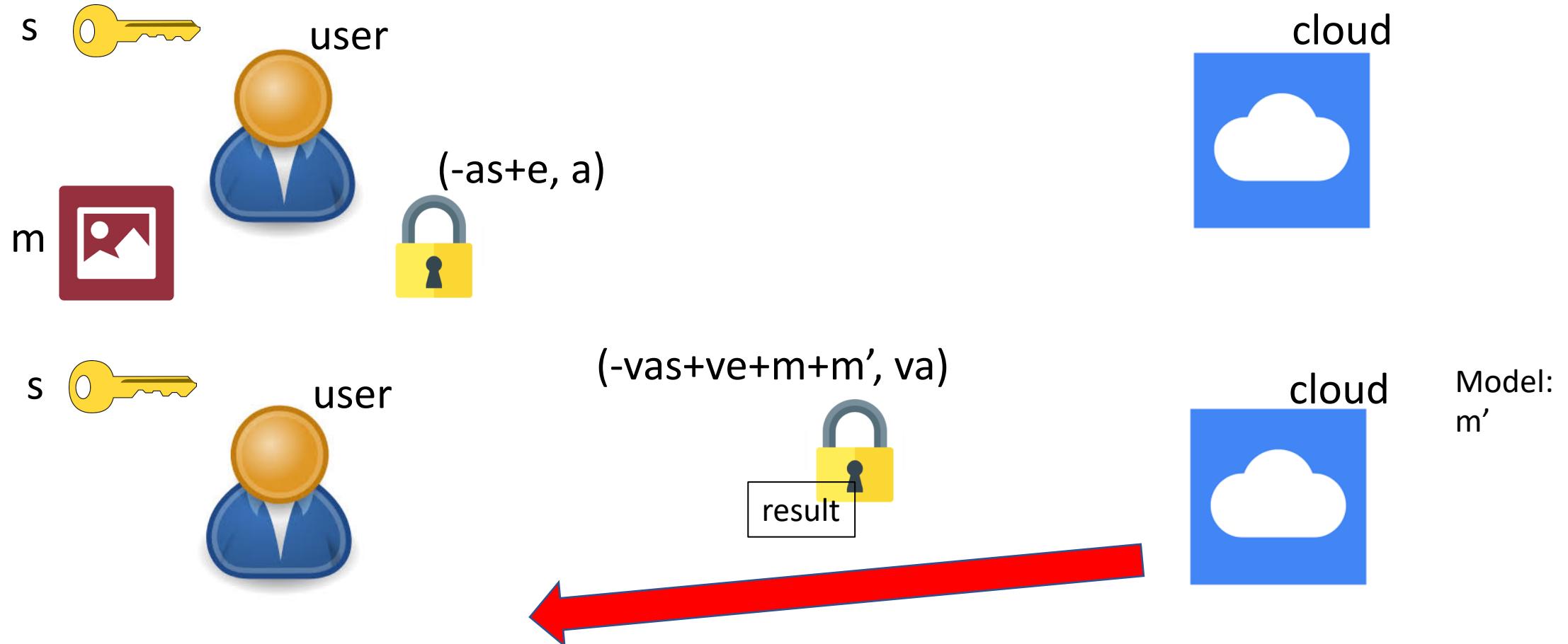
HE scenario 1 (Cloud inference)



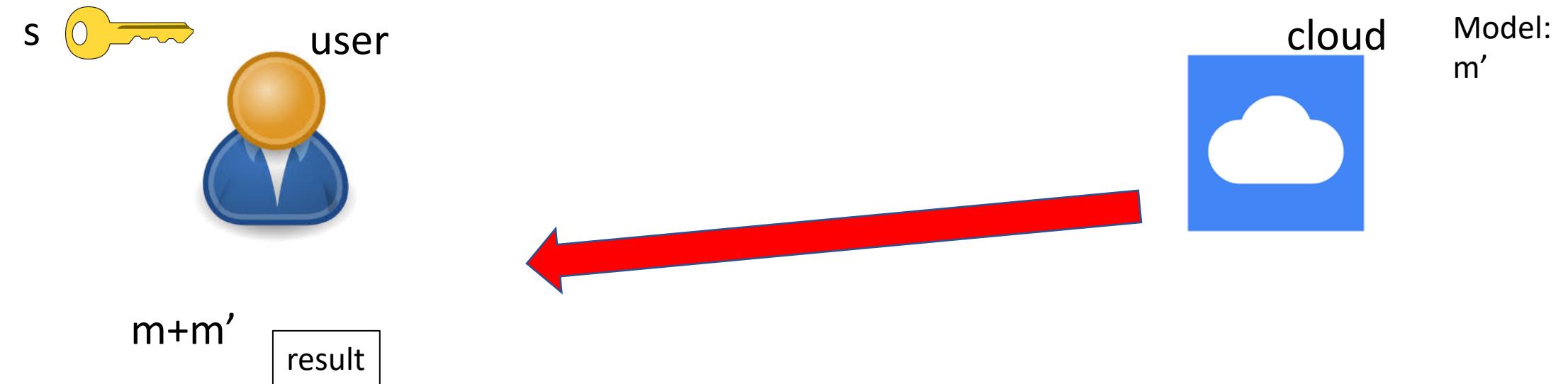
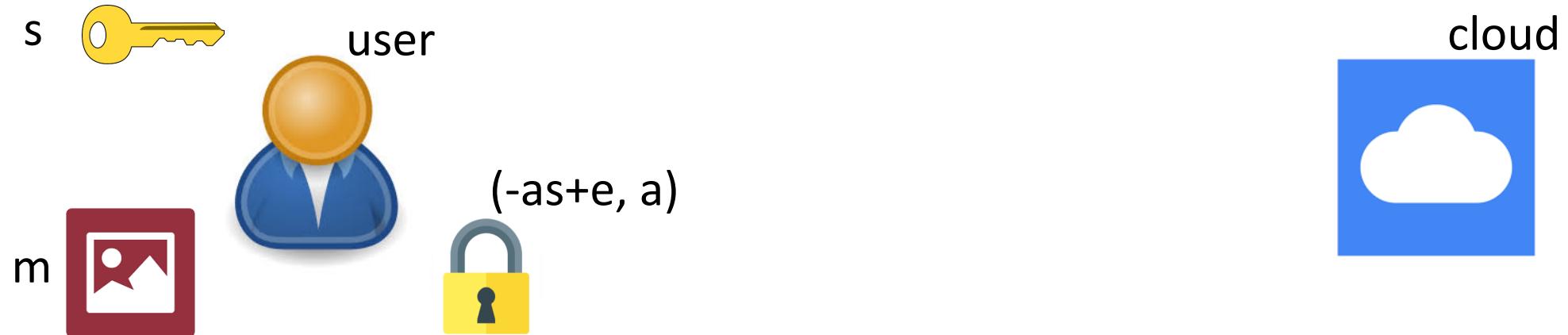
HE scenario 1 (Cloud inference)



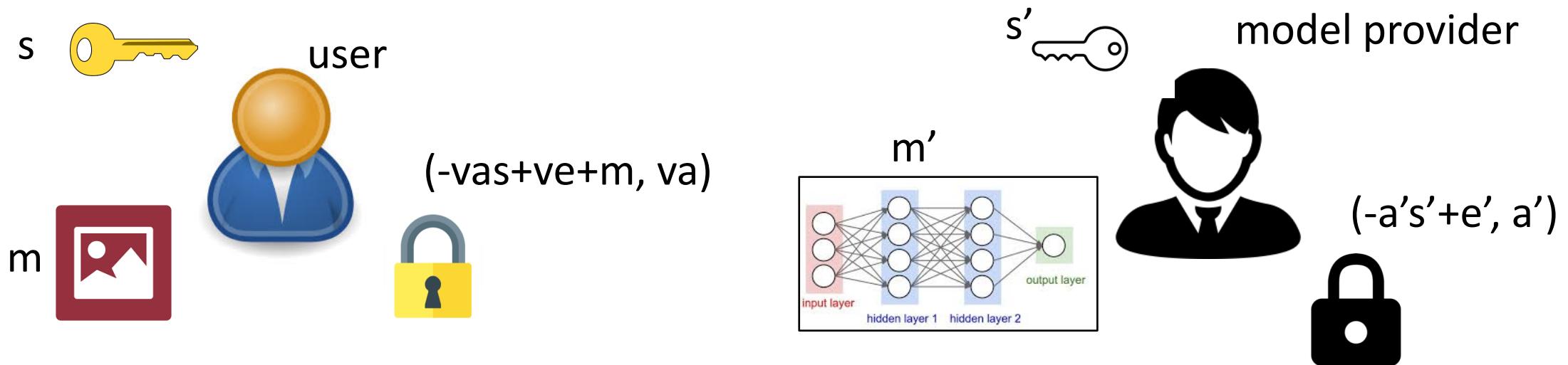
HE scenario 1 (Cloud inference)



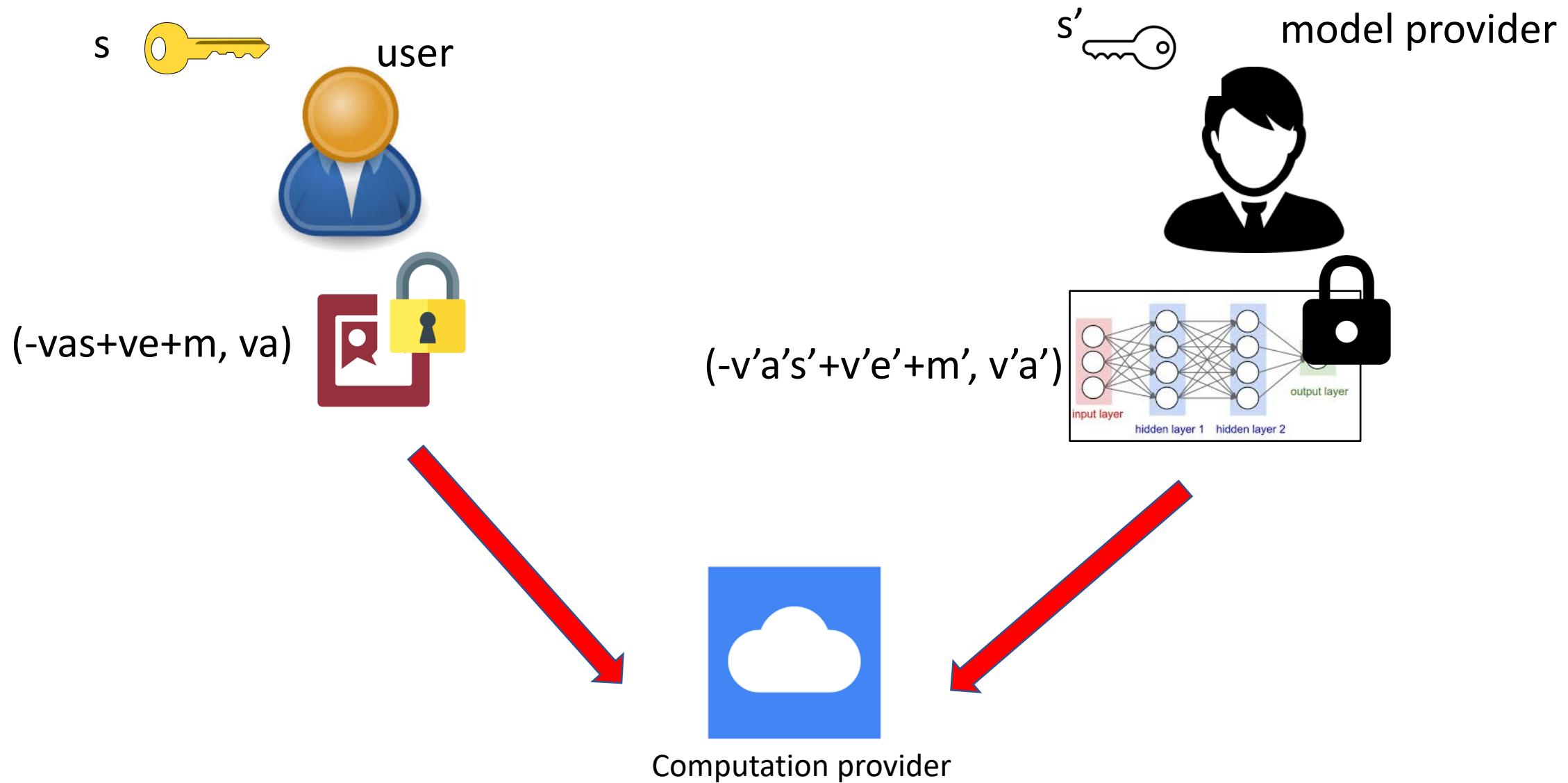
HE scenario 1 (Cloud inference)



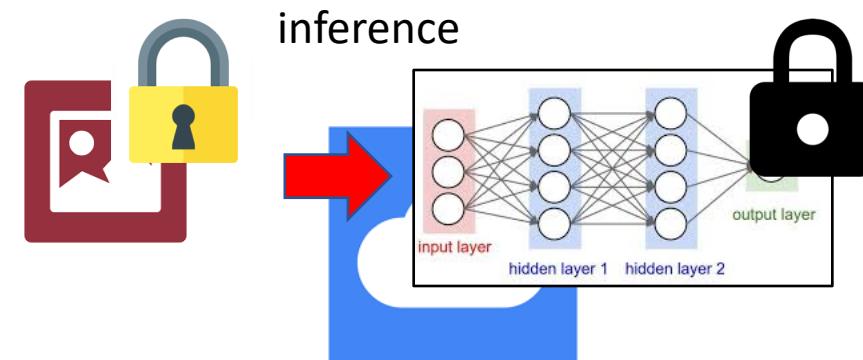
HE scenario 2 (Cloud computation)



HE scenario 2 (Cloud computation)



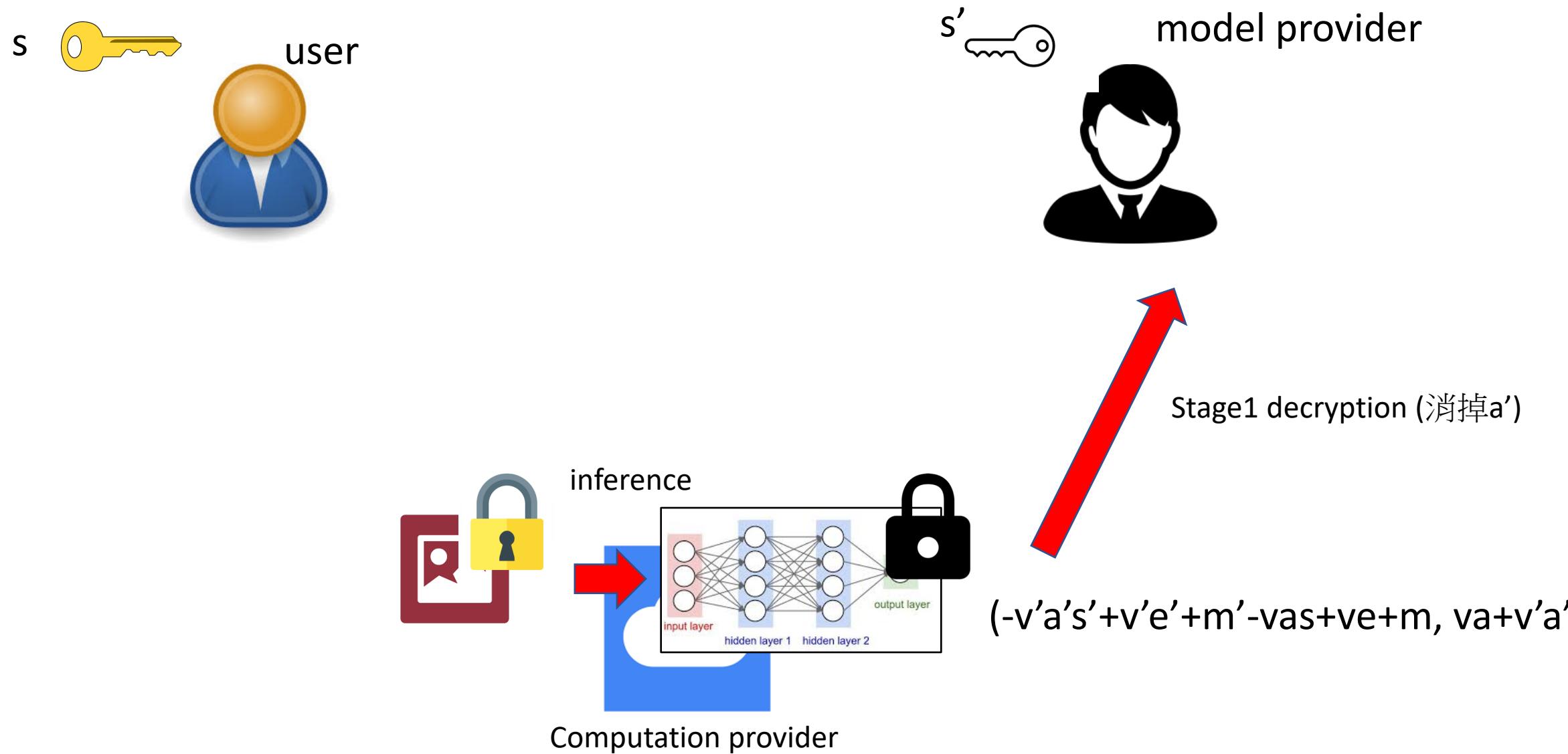
HE scenario 2 (Cloud computation)



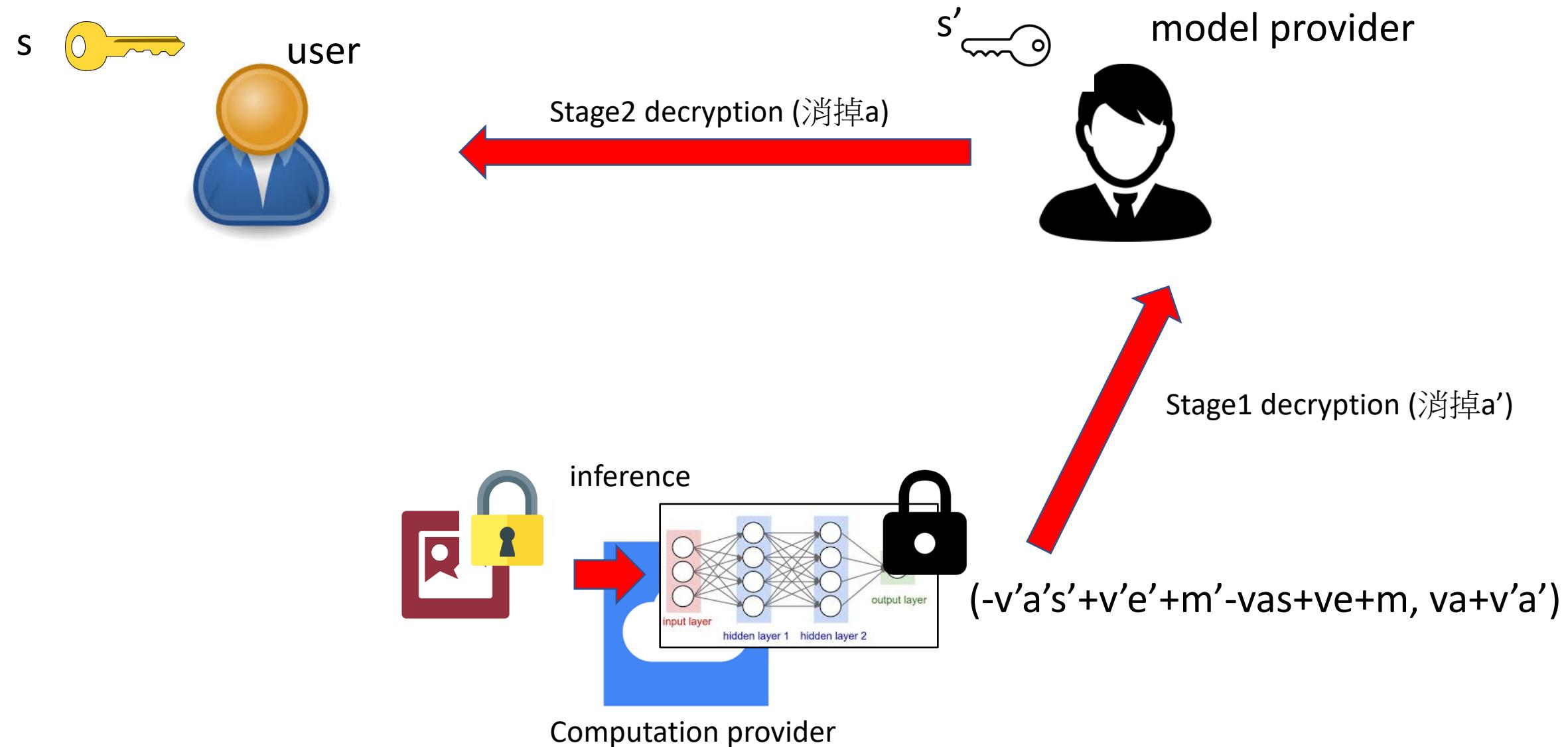
Computation provider

$$(-v'a's' + v'e' + m' - vas + ve + m, va + v'a')$$

HE scenario 2 (Cloud computation)



HE scenario 2 (Cloud computation)



Possible applications

Domain	Genomics	Health	National Security	Education
Topic	Match Maker	Billing and Reporting	Smart Grid (Municipal Service)	School Dropouts
Data Owner	Medical Institutions	Small Hospital, Clinic	Nodes and Network	School, Hospital, Welfare
Latency of Service	Hours	Hours	Quasi-Real Time	Week
Data volume (size x no)	DB O(1000X1MB) Query O(1KB)	O(10M) x O(1M)	O(1M) x O(1M)	O(1K) x O(1M)
Data persistency	Add Only	Add Only	Add Only	Add Only
Technical Issues	Comparison, Sorting Auditing Privacy	Tabulation, Linear Algebra	Comparison	Comparison Matrix Analysis
When is possible	1 years	2-3 years	Now	2-3 years
Why HE?	HIPAA	Cyber insurance	Privacy	FERPA
Who pays?	Health Insurance	Hospital	Energy Company	DoE

https://dualitytech.com/wp-content/uploads/2020/04/Covid19_onepager.pdf

https://www.researchgate.net/publication/320976976_APPLICATIONS_OF_HOMOMORPHIC_ENCRYPTION

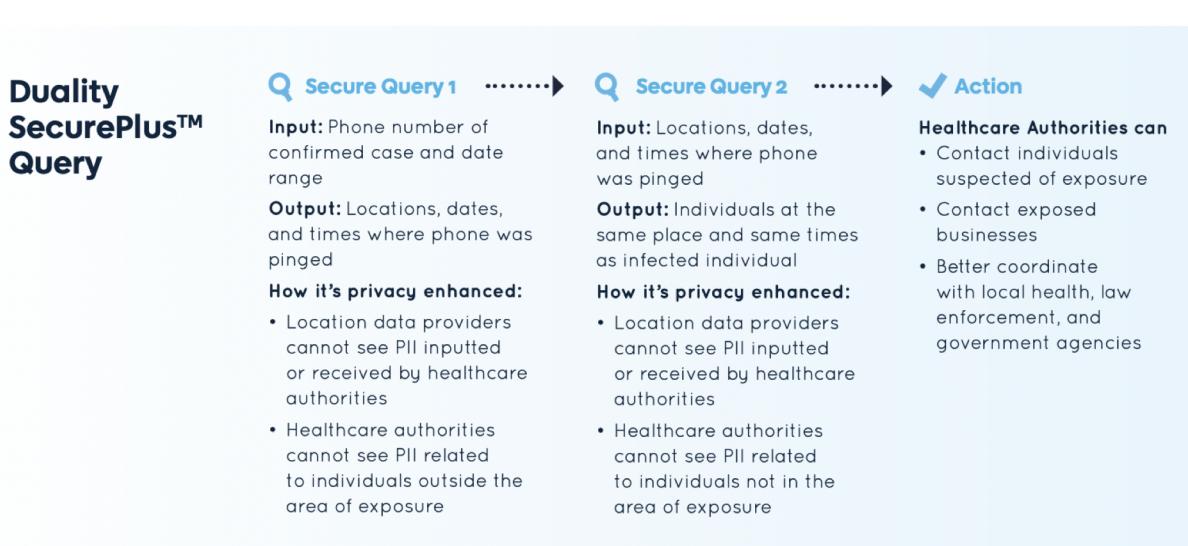
Possible applications

Domain	Genomics	Health	National Security	Education
Topic	Match Maker	Billing and Reporting	Smart Grid (Municipal Service)	School Dropouts
Data Owner	Medical Institutions	Small Hospital, Clinic	Nodes and Network	School, Hospital, Welfare
Latency of Service	Hours	Hours	Quasi-Real Time	Week
Data volume (size x no)	DB O(1000X1MB) Query O(1KB)	O(10M) x O(1M)	O(1M) x O(1M)	O(1K) x O(1M)
Data persistency	Add Only	Add Only	Add Only	Add Only
Technical Issues	Comparison, Sorting Auditing Privacy	Tabulation, Linear Algebra	Comparison	Comparison Matrix Analysis
When is possible	1 years	2-3 years	Now	2-3 years
Why HE?	HIPAA	Cyber insurance	Privacy	FERPA
Who pays?	Health Insurance	Hospital	Energy Company	DoE

https://dualitytech.com/wp-content/uploads/2020/04/Covid19_onepager.pdf

https://www.researchgate.net/publication/320976976_APPLICATIONS_OF_HOMOMORPHIC_ENCRYPTION

Privacy-Preserving, Large-Scale COVID-19 Contact Tracing



Partners and Recognition



Possible applications

Duality SecurePlus™ Query

感染者的電話號碼和日期範圍(加密的)
→感染者什麼時候經過哪裡

🔍 Secure Query 1 ➔

Input: Phone number of confirmed case and date range

Output: Locations, dates, and times where phone was pinged

How it's privacy enhanced:

- Location data providers cannot see PII inputted or received by healthcare authorities
- Healthcare authorities cannot see PII related to individuals outside the area of exposure

時間地點(加密的)
→哪些人在同時間經過這裡

🔍 Secure Query 2 ➔

Input: Locations, dates, and times where phone was pinged

Output: Individuals at the same place and same times as infected individual

How it's privacy enhanced:

- Location data providers cannot see PII inputted or received by healthcare authorities
- Healthcare authorities cannot see PII related to individuals not in the area of exposure

✓ Action

Healthcare Authorities can

- Contact individuals suspected of exposure
- Contact exposed businesses
- Better coordinate with local health, law enforcement, and government agencies

PII: Personal Identifiable Information

https://dualitytech.com/wp-content/uploads/2020/04/Covid19_onepager.pdf

https://www.researchgate.net/publication/320976976_APPLICATIONS_OF_HOMOMORPHIC_ENCRYPTION

Open source

Name	Developer	BGV ^[17]	CKKS ^[35]	BFV ^[20]	FHEW ^[32]	CKKS Bootstrapping ^[39]	TFHE ^[33]	Description
HElib ^[40]	IBM	Yes	Yes	No	No	No	No	BGV scheme with the GHS optimizations.
Microsoft SEAL ^[41]	Microsoft	No	Yes	Yes	No	No	No	
PALISADE ^[42]	Consortium of DARPA–funded defense contractors and academics: New Jersey Institute of Technology, Duality Technologies, Raytheon BBN Technologies, MIT, University of California, San Diego and others.	Yes	Yes	Yes	Yes	No	Yes	General-purpose lattice cryptography library.
HEAAN ^[43]	Seoul National University	No	Yes	No	No	Yes	No	
FHEW ^[32]	Leo Ducas and Daniele Micciancio	No	No	No	Yes	No	No	
TFHE ^[33]	Ilaria Chillotti, Nicolas Gama, Mariya Georgieva and Malika Izabachene	No	No	No	No	No	Yes	
FV-NFLlib ^[44]	CryptoExperts	No	No	Yes	No	No	No	
NuFHE ^[45]	NuCypher	No	No	No	No	No	Yes	Provides a GPU implementation of TFHE.
Lattigo ^[46]		No	Yes	Yes	No	No	No	Implementation in Go along with their distributed variants enabling Secure multi–party computation.

Homomorphic encryption summary

- Privacy preserving computation
 - Many use cases
- Huge computational cost (10000x)
 - GPU
 - FPGA
 - ASIC

Outline

- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

RSA encrypt

- $c = Enc(m) = m^e \text{ mod } n, n = pq$
- $Enc(m_1) * Enc(m_2) = (m_1 m_2)^e = Enc(m_1 * m_2)$
- Example
 - $c = Enc(m) = m^5 \text{ mod } 35$
 - $m_1 = 2, m_2 = 3$
 - $Enc(m_1) = 32, Enc(m_2) = 33$
 - $Enc(m_1) * Enc(m_2) = 16 * 11 \text{ mod } 35 = 6 = 2 * 3$
- Only multiplication

Lattice-based cryptography

- Multiplication and addition
- Quantum-resistant
- Fair enough speed and memory usage
- Easy to understand

Computation on $\frac{\mathbb{Z}_q}{x^n + 1}$

- Example $\frac{\mathbb{Z}_7}{x^4 + 1}$
 - $a \rightarrow 6x^3 + 4x^2 + 1x^1 + 2x^0$
 - $b \rightarrow 1x^3 + 1x^2 + 0x^1 + 1x^0$
 - $a + b \rightarrow 0x^3 + 5x^2 + 1x^1 + 3x^0$

Computation on $\frac{\mathbb{Z}_q}{x^n+1}$

- Example $\frac{\mathbb{Z}_7}{x^4+1}$

- $a \rightarrow 6x^3 + 4x^2 + 1x^1 + 2x^0$
- $b \rightarrow 1x^3 + 1x^2 + 0x^1 + 1x^0$
- $a + b \rightarrow 0x^3 + 5x^2 + 1x^1 + 3x^0$
- $a * b \rightarrow 6x^6 + 10x^5 + 5x^4 + 9x^3 + 6x^2 + 1x^1 + 2x^0$
 $\quad \rightarrow 6(-x^2) + 10(-x) + 5(-1) + 9x^3 + 6x^2 + 1x^1 + 2x^0$
 $\quad \rightarrow 9x^3 + 0x^2 - 9x^1 - 3x^0 \text{ (plug in } x^4 = -1\text{)}$
 $\quad \rightarrow 2x^3 + 0x^2 + 5x^1 + 4x^0 \pmod{7}$

Computation on $\frac{\mathbb{Z}_q}{x^n+1}$

- Example $\frac{\mathbb{Z}_7}{x^4+1}$

- $a \rightarrow 6x^3 + 4x^2 + 1x^1 + 2x^0$
- $b \rightarrow 1x^3 + 1x^2 + 0x^1 + 1x^0$
- $a + b \rightarrow 0x^3 + 5x^2 + 1x^1 + 3x^0$
- $a * b \rightarrow 6x^6 + 10x^5 + 5x^4 + 9x^3 + 6x^2 + 1x^1 + 2x^0$
 $\rightarrow 6(-x^2) + 10(-x) + 5(-1) + 9x^3 + 6x^2 + 1x^1 + 2x^0$
 $\rightarrow 9x^3 + 0x^2 - 9x^1 - 3x^0 \text{ (plug in } x^4 = -1\text{)}$
 $\rightarrow 2x^3 + 0x^2 + 5x^1 + 4x^0 \pmod{7}$

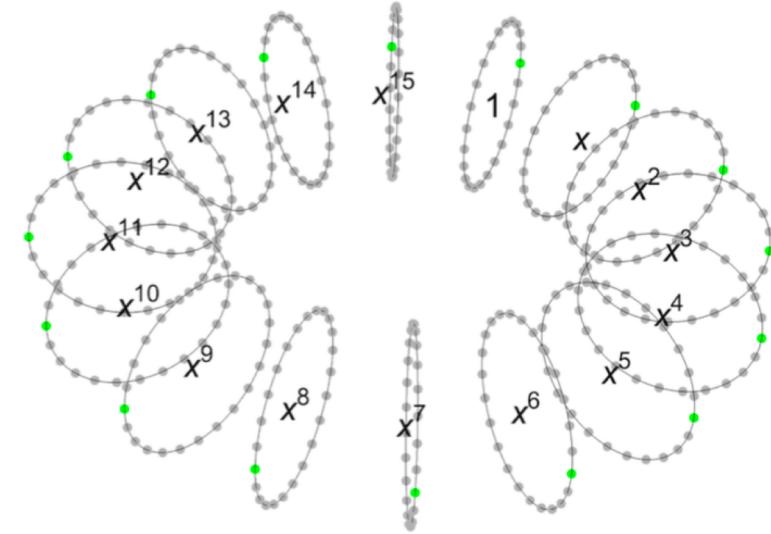


Figure 4.1: Depiction of a polynomial ring with degree 16 and plaintext modulus 24 [35].

BFV Keygen

- Plaintext domain in $\frac{Z_t}{x^n + 1}$, encrypted domain in $\frac{Z_q}{x^n + 1}$, $q \gg t$

BFV Keygen

- Plaintext domain in $\frac{Z_t}{x^n+1}$, encrypted domain in $\frac{Z_q}{x^n+1}$, $q \gg t$
- Example $\frac{Z_7}{x^4+1}$ ($t=7$, 3bits), $\frac{Z_{1153}}{x^4+1}$ ($q=1153$, 10bits), $\frac{q}{t} = 164$ (7 bits)

BFV Keygen

- Plaintext domain in $\frac{\mathbb{Z}_t}{x^n+1}$, encrypted domain in $\frac{\mathbb{Z}_q}{x^n+1}$, $q \gg t$
- Example $\frac{\mathbb{Z}_7}{x^4+1}$ ($t=7$, 3bits), $\frac{\mathbb{Z}_{1153}}{x^4+1}$ ($q=1153$, 10bits), $\frac{q}{t} = 164$ (7 bits)
- Keygen
 - $a = 123x^3 + 456x^2 + 789x^1 + 12x^0 \in \frac{\mathbb{Z}_{1153}}{x^4+1}$ (encrypted domain)
 - $s = 0x^3 + 1x^2 + 0x^1 + 1x^0$ (**secret key**, binary coefficients)
 - $e = 1x^3 + 0x^2 + 1x^1 + 0x^0$ (**noise**, binary coefficients)
 - $pk = (e - a * s, a)$. ($a * s = 912x^3 + 468x^2 + 666x^1 + 709x^0$)

BFV Keygen

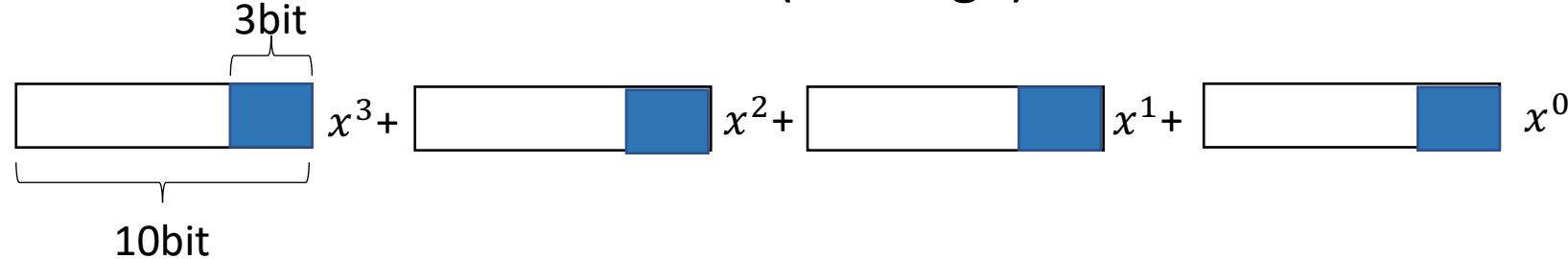
- Plaintext domain in $\frac{\mathbb{Z}_t}{x^n+1}$, encrypted domain in $\frac{\mathbb{Z}_q}{x^n+1}$, $q \gg t$
- Example $\frac{\mathbb{Z}_7}{x^4+1}$ ($t=7$, 3bits), $\frac{\mathbb{Z}_{1153}}{x^4+1}$ ($q=1153$, 10bits), $\frac{q}{t} = 164$ (7 bits)
- Keygen
 - $a = 123x^3 + 456x^2 + 789x^1 + 12x^0 \in \frac{\mathbb{Z}_{1153}}{x^4+1}$ (encrypted domain)
 - $s = 0x^3 + 1x^2 + 0x^1 + 1x^0$ (**secret key**, binary coefficients)
 - $e = 1x^3 + 0x^2 + 1x^1 + 0x^0$ (**noise**, binary coefficients)
 - $pk = (e - a * s, a)$. ($a * s = 912x^3 + 468x^2 + 666x^1 + 709x^0$)
 - Only the party with s can cancel out the term $a * s$

BFV Keygen

- Plaintext domain in $\frac{\mathbb{Z}_t}{x^n+1}$, encrypted domain in $\frac{\mathbb{Z}_q}{x^n+1}$, $q \gg t$
- Example $\frac{\mathbb{Z}_7}{x^4+1}$ ($t=7$, 3bits), $\frac{\mathbb{Z}_{1153}}{x^4+1}$ ($q=1153$, 10bits), $\frac{q}{t} = 164$ (7 bits)
- Keygen
 - $a = 123x^3 + 456x^2 + 789x^1 + 12x^0 \in \frac{\mathbb{Z}_{1153}}{x^4+1}$ (encrypted domain)
 - $s = 0x^3 + 1x^2 + 0x^1 + 1x^0$ (**secret key**, binary coefficients)
 - $e = 1x^3 + 0x^2 + 1x^1 + 0x^0$ (**noise**, binary coefficients)
 - $pk = (e - a * s, a)$. ($a * s = 912x^3 + 468x^2 + 666x^1 + 709x^0$)
 - Only the party with s can cancel out the term $a * s$
 - Hardness: given pk , hard to find s (reduce to SVP in lattice which is NP-Hard)

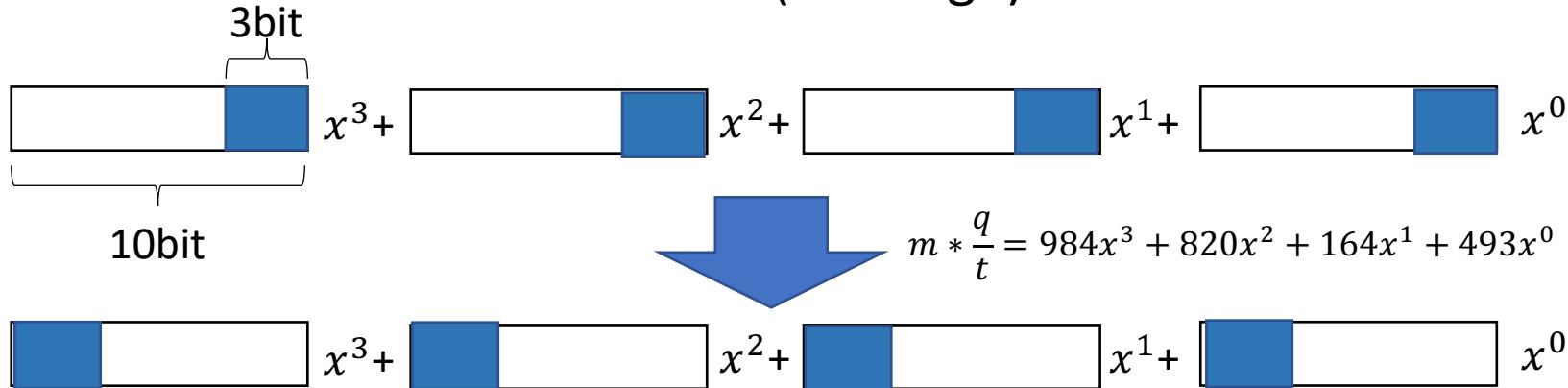
BFV Encrypt

- Encrypt $(m * \frac{q}{t} + pk[0], pk[1]) = (m * \frac{q}{t} + (e - a * s), a)$
 - $m = 6x^3 + 5x^2 + 1x^1 + 3x^0$ (message)



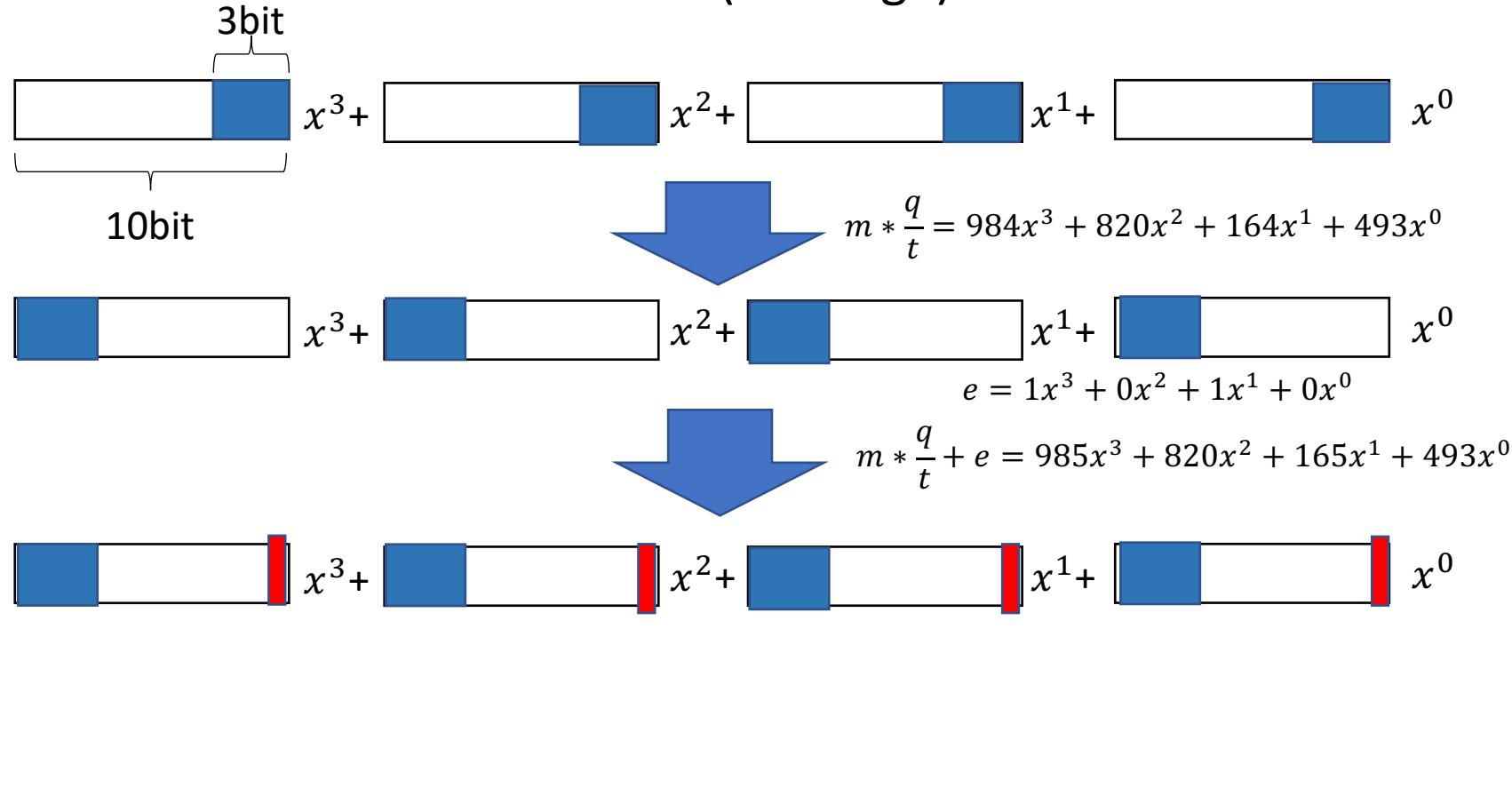
BFV Encrypt

- Encrypt $(m * \frac{q}{t} + pk[0], pk[1]) = (m * \frac{q}{t} + (e - a * s), a)$
 - $m = 6x^3 + 5x^2 + 1x^1 + 3x^0$ (message)



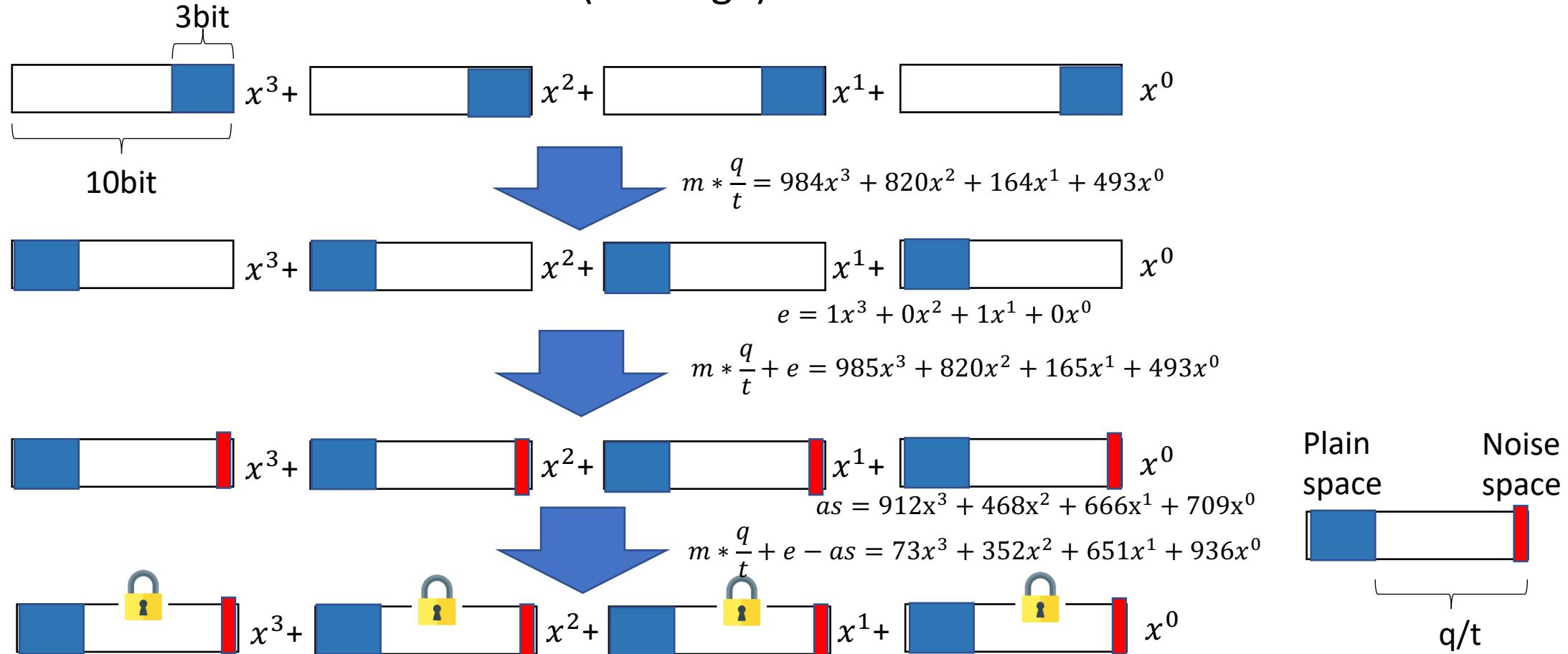
BFV Encrypt

- Encrypt $(m * \frac{q}{t} + pk[0], pk[1]) = (m * \frac{q}{t} + (e - a * s), a)$
- $m = 6x^3 + 5x^2 + 1x^1 + 3x^0$ (message)



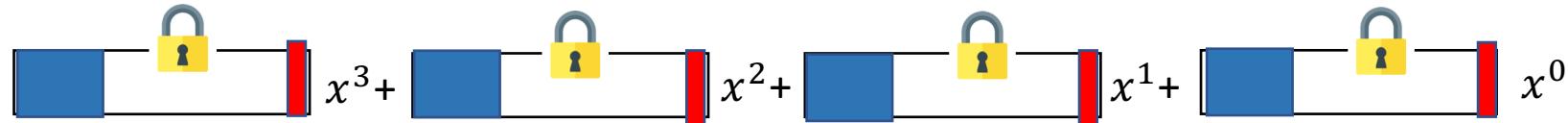
BFV Encrypt

- Encrypt $(m * \frac{q}{t} + pk[0], pk[1]) = (m * \frac{q}{t} + (e - a * s), a)$
- $m = 6x^3 + 5x^2 + 1x^1 + 3x^0$ (message)



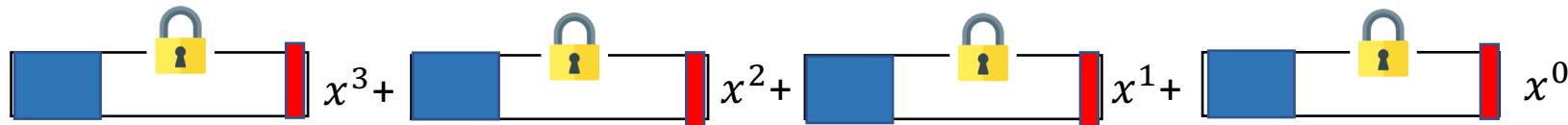
BFV Decrypt

$$m * \frac{q}{t} + e - as = 73x^3 + 352x^2 + 651x^1 + 936x^0$$

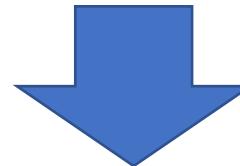


BFV Decrypt

$$m * \frac{q}{t} + e - as = 73x^3 + 352x^2 + 651x^1 + 936x^0$$



$$as = 912x^3 + 468x^2 + 666x^1 + 709x^0$$

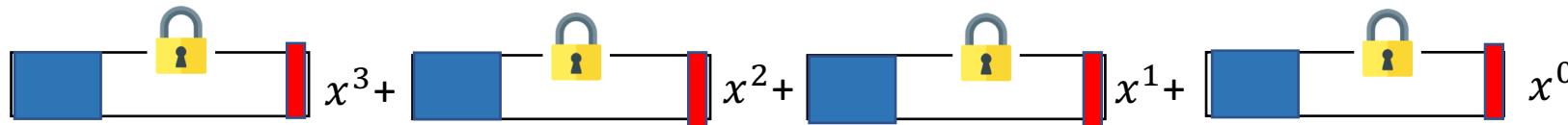


$$m * \frac{q}{t} + e - a * s + a * s = 985x^3 + 820x^2 + 164x + 492$$



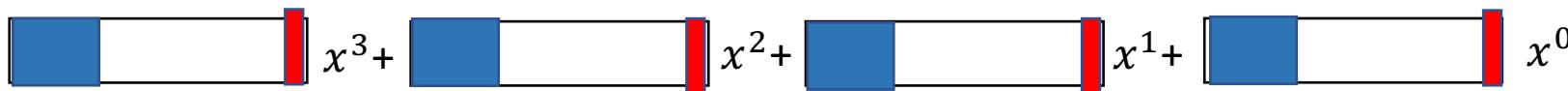
BFV Decrypt

$$m * \frac{q}{t} + e - as = 73x^3 + 352x^2 + 651x^1 + 936x^0$$



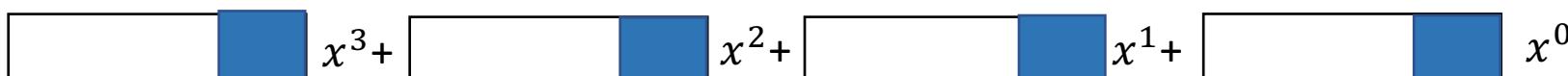
$$as = 912x^3 + 468x^2 + 666x^1 + 709x^0$$

$$m * \frac{q}{t} + e - a * s + a * s = 985x^3 + 820x^2 + 164x + 492$$



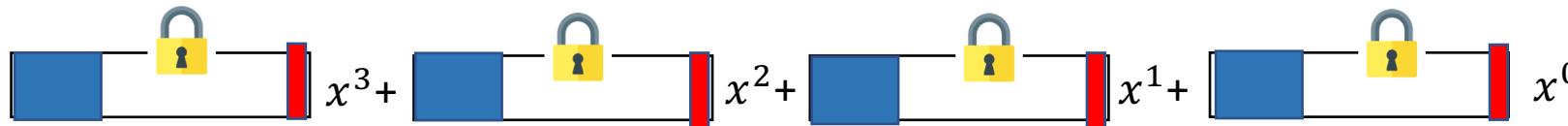
$$t/q = 0.00607111882$$

$$\left(m * \frac{q}{t} + e \right) * \frac{t}{q} = 5.98005203769999x^3 + 4.9783174324x^2 + 0.9956634864799999x^1 + 2.98699045944$$



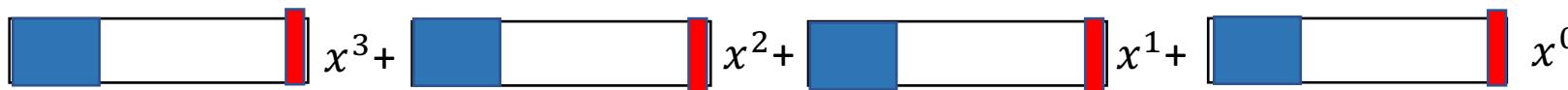
BFV Decrypt

$$m * \frac{q}{t} + e - as = 73x^3 + 352x^2 + 651x^1 + 936x^0$$



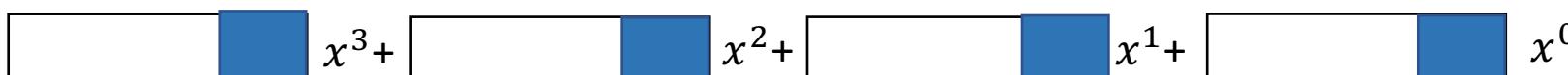
$$as = 912x^3 + 468x^2 + 666x^1 + 709x^0$$

$$m * \frac{q}{t} + e - a * s + a * s = 985x^3 + 820x^2 + 164x + 492$$



$$t/q=0.00607111882$$

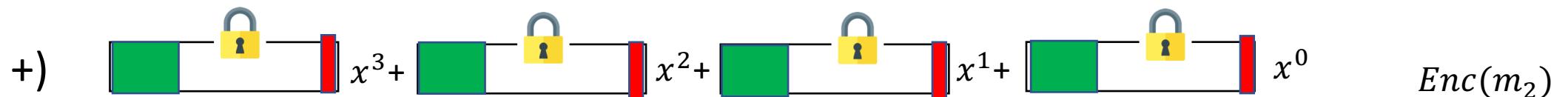
$$\left(m * \frac{q}{t} + e \right) * \frac{t}{q} = 5.98005203769999x^3 + 4.9783174324x^2 + 0.9956634864799999x^1 + 2.98699045944$$



$$\text{After round: } 6x^3 + 5x^2 + 1x + 3$$

Decryption succeed if $e < \frac{q}{2t}$

Homomorphic addition



Decryption succeed if $e < \frac{q}{2t}$

Plaintext addition



+)

The diagram shows the encryption of message m_2 as a polynomial over $\text{GF}(2)$. The polynomial is represented as $x^3 + x^2 + x^1 + x^0$. The first term is green, while the others are white rectangles. To the right of the polynomial is the label m_2 .



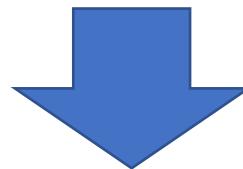
Decryption succeed if $e < \frac{q}{2t}$

Plaintext multiplication



+)

	green	$x^3 +$	green	$x^2 +$	green	$x^1 +$	green	x^0	m_2
--	-------	---------	-------	---------	-------	---------	-------	-------	-------



Decryption succeed if $e < \frac{q}{2t}$

Homomorphic multiplication

- Skip, a little bit complicated
- **Consume lots of noise space**

Homomorphic multiplication

- Skip, a little bit complicated
 - **Consume lots of noise space**

t = 1153 (10bits)

$q = 1427247692705959881058285969449495136382747009$ (150bits)

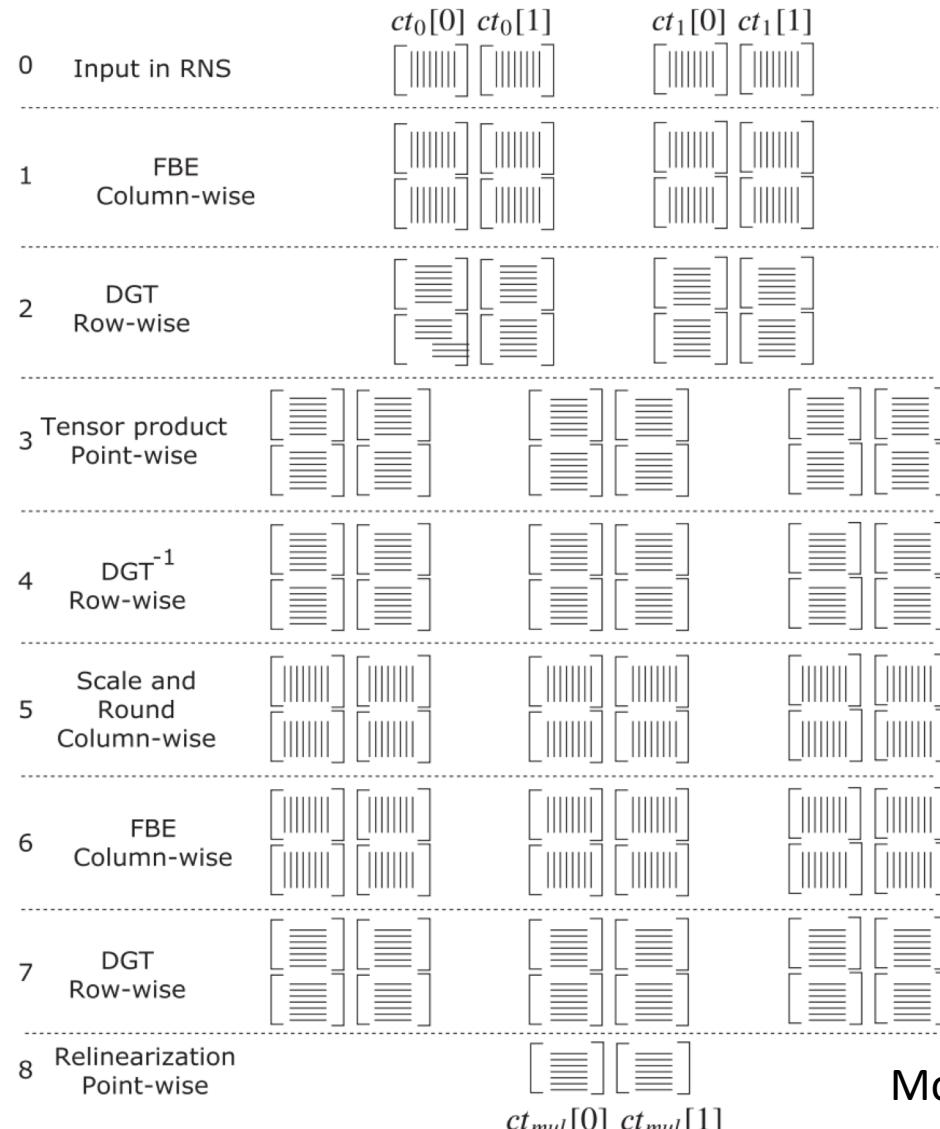
n = 16

Before round, apply HMul several times:

x^0 coefficient

Noise bits

Homomorphic multiplication

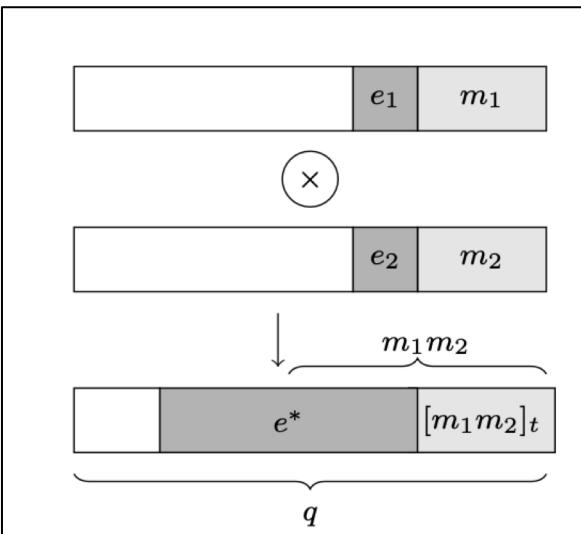


Normal NTT vs DGT (GPU-friendly)

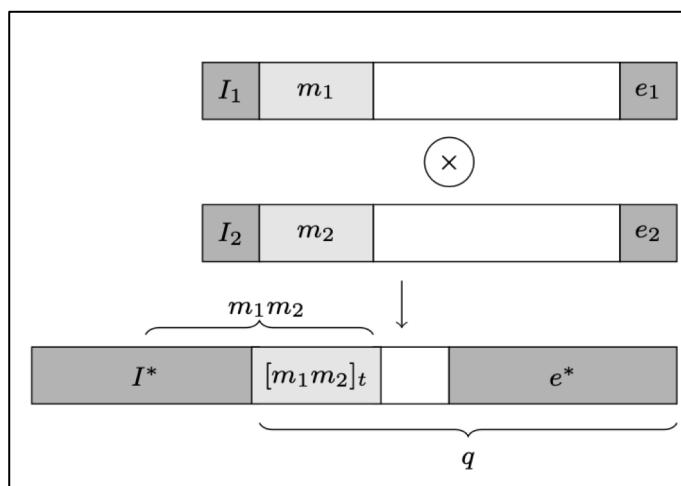
Modulus raising vs bit-decomposition (FPGA-friendly)

Comparison between different schemes

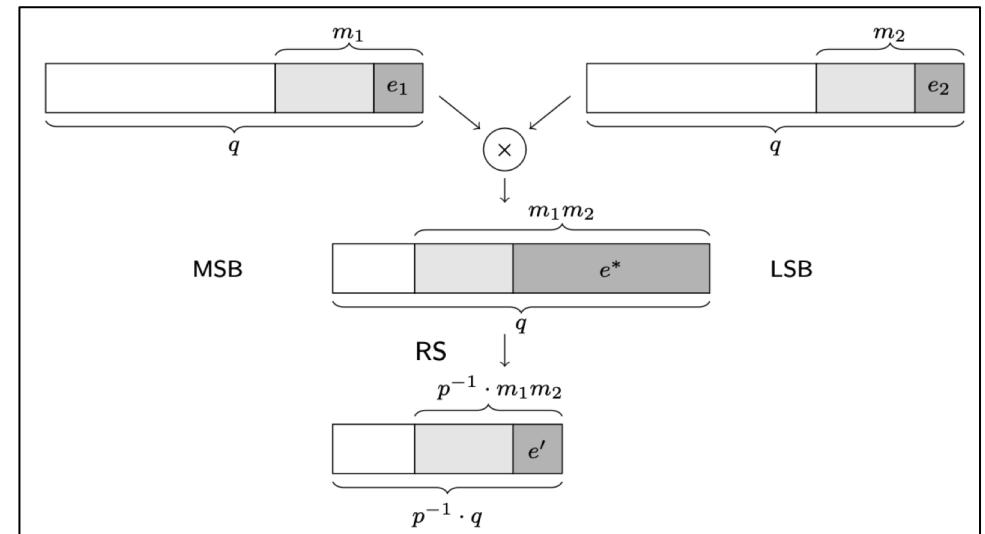
- Decryption structure



BGV, MSB is truncated



BFV, MSB is truncated
(suitable for integer operation,
can also perform fixed point
operation using $m * 2^{\{precision\}}$)



CKKS, LSB is truncated (suitable for **approximate computation**)

Encoding

$$\frac{Z_{17}}{x^4 + 1}$$

$$m_1 = 6x^3 + 5x^2 + 1x^1 + 3x^0$$

$$m_2 = 4x^3 + 5x^2 + 1x^1 + 2x^0$$

$$m_1 m_2 = 0x^3 + 2x^2 + 6x^1 + 5x^0$$

Not that useful

Encoding

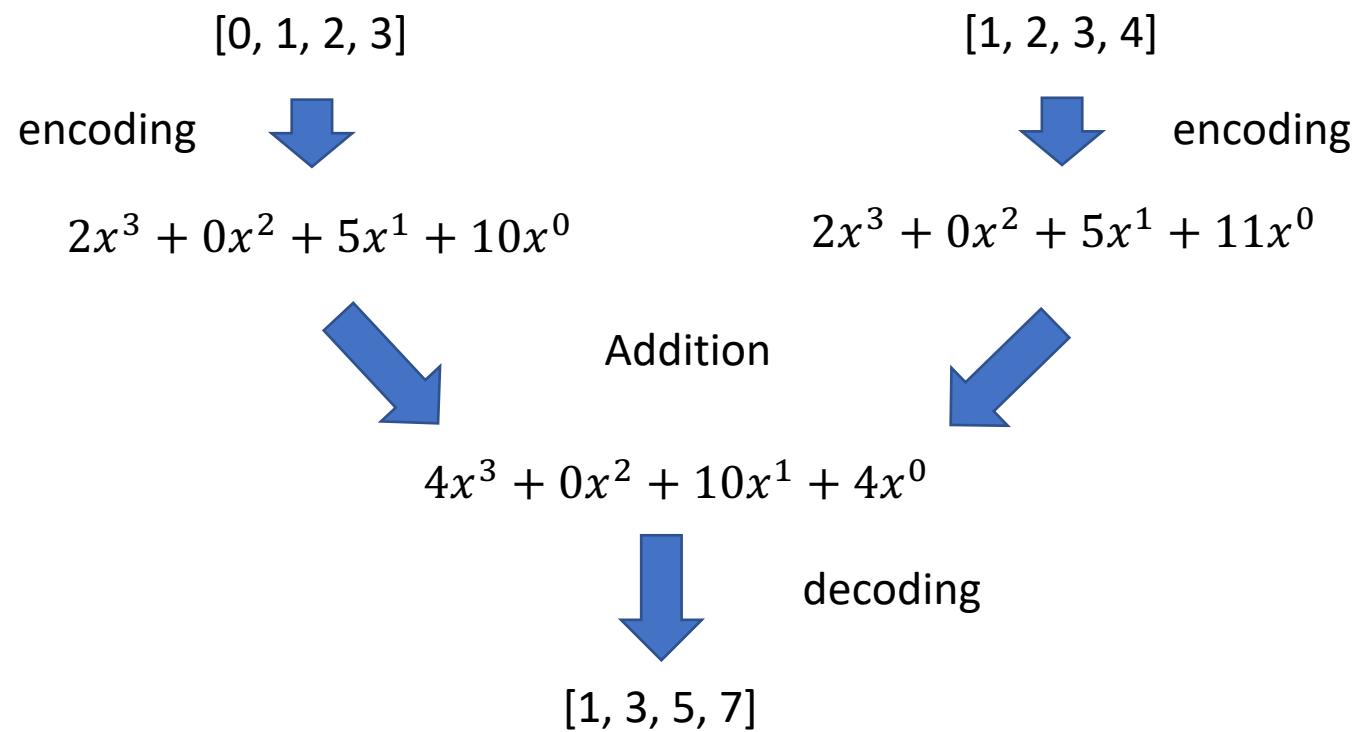
$$\frac{Z_{17}}{x^4 + 1}$$

$$m_1 = 6x^3 + 5x^2 + 1x^1 + 3x^0$$

$$m_2 = 4x^3 + 5x^2 + 1x^1 + 2x^0$$

$$m_1 m_2 = 0x^3 + 2x^2 + 6x^1 + 5x^0$$

Not that useful



Encoding

Require $t \equiv 1 \pmod{2n}$

Encode

$$13 \begin{pmatrix} 9^{-1*0} & 9^{-3*0} & 9^{-7*0} & 9^{-5*0} \\ 9^{-1*1} & 9^{-3*1} & 9^{-7*1} & 9^{-5*1} \\ 9^{-1*2} & 9^{-3*2} & 9^{-7*2} & 9^{-5*2} \\ 9^{-1*3} & 9^{-3*3} & 9^{-7*3} & 9^{-5*3} \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 2 \\ 3 \end{pmatrix} = \begin{pmatrix} 10 \\ 5 \\ 0 \\ 2 \end{pmatrix}$$

Decode

$$\begin{pmatrix} 9^{1*0} & 9^{1*1} & 9^{1*2} & 9^{1*3} \\ 9^{3*0} & 9^{3*1} & 9^{3*2} & 9^{3*3} \\ 9^{7*0} & 9^{7*1} & 9^{7*2} & 9^{7*3} \\ 9^{5*0} & 9^{5*1} & 9^{5*2} & 9^{5*3} \end{pmatrix} \begin{pmatrix} 4 \\ 10 \\ 0 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 \\ 3 \\ 5 \\ 7 \end{pmatrix}$$

$$\frac{Z_{17}}{x^4 + 1}$$

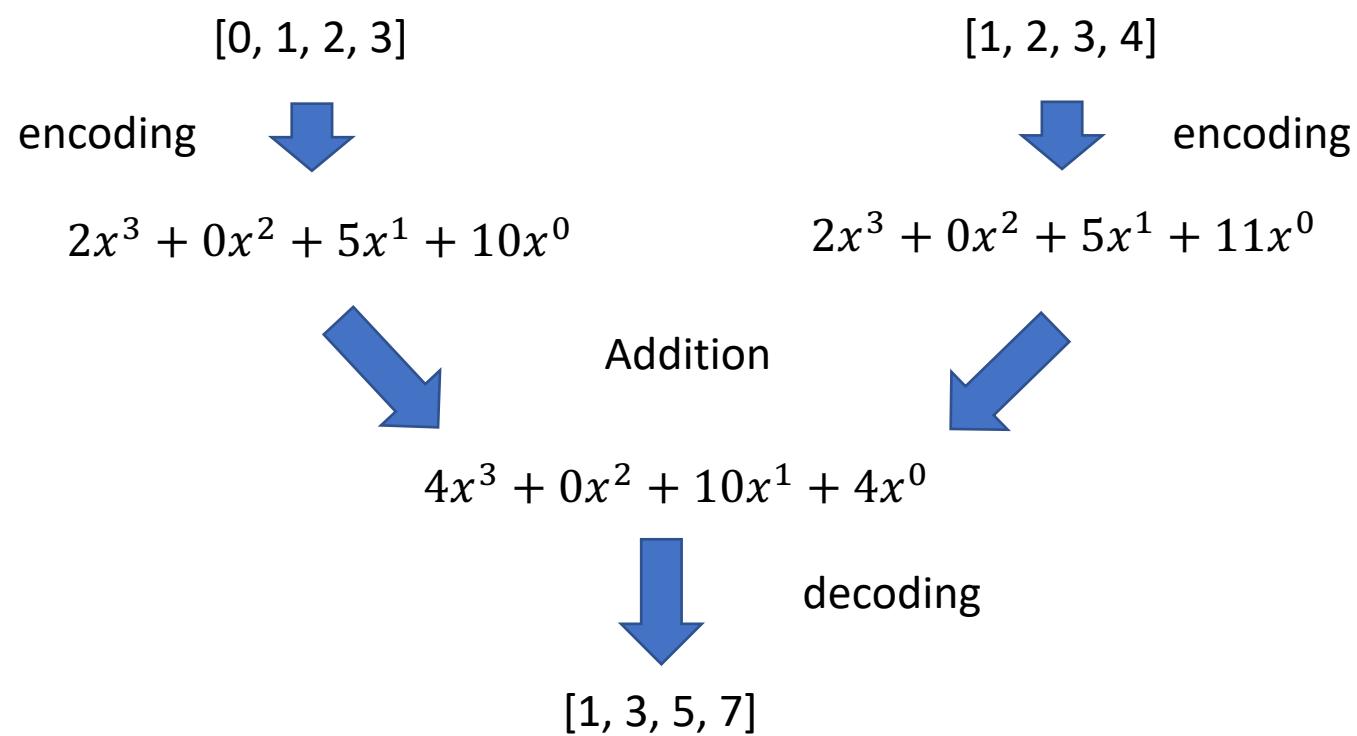
Require $t \equiv 1 \pmod{2n}$

$$m_1 = 6x^3 + 5x^2 + 1x^1 + 3x^0$$

$$m_2 = 4x^3 + 5x^2 + 1x^1 + 2x^0$$

$$m_1 m_2 = 0x^3 + 2x^2 + 6x^1 + 5x^0$$

Not that useful



Rotation

$$\frac{Z_{17}}{x^8 + 1}$$

[0, 1, 2, 3]
[4, 5, 6, 7]

Rotate column



[3, 0, 1, 2]
[7, 4, 5, 6]

$$\frac{Z_{17}}{x^8 + 1}$$

[0, 1, 2, 3]
[4, 5, 6, 7]

Rotate row



[4, 5, 6, 7]
[0, 1, 2, 3]

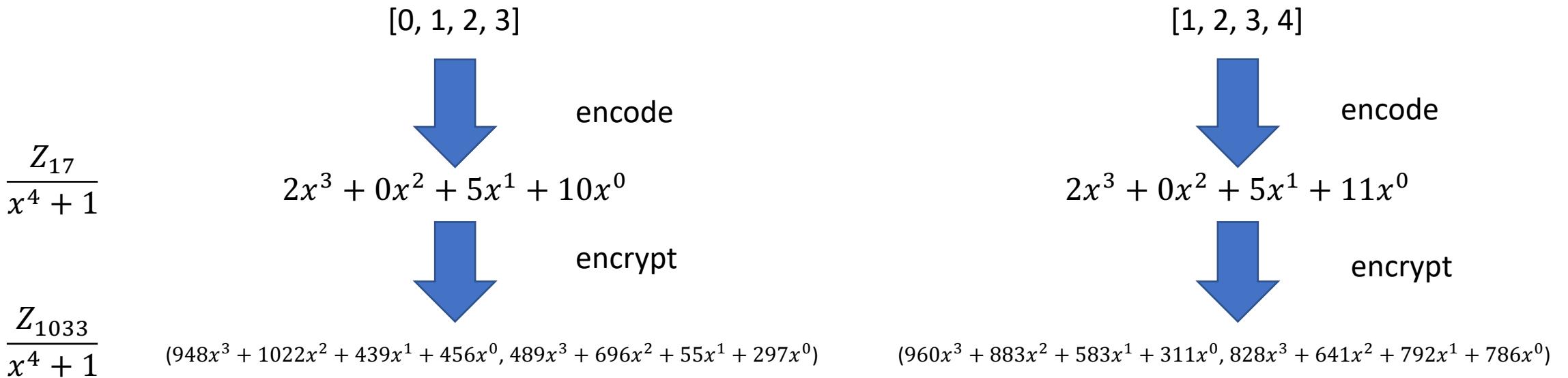
Rotation

- Using a set of automorphism R_{n^r} of the form

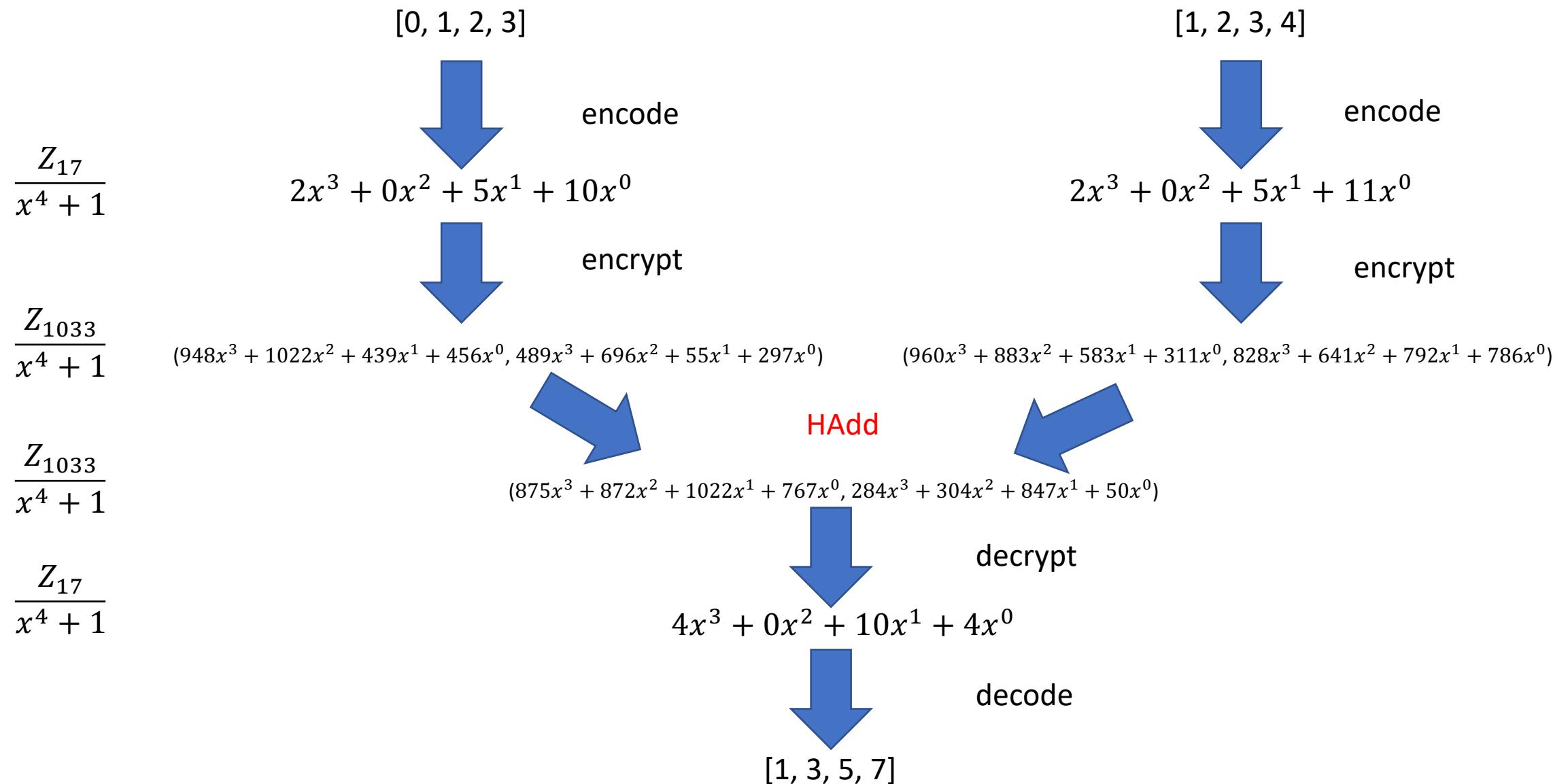
$$\theta_t : R_{p^r} \longrightarrow R_{p^r}, \quad a(X) \longmapsto a(X^t) \pmod{(p^r, \Phi_m(X))}.$$

- Example: $\frac{Z_{17}}{x^4+1}$
 - $2x^3 + 3x^2 + 1x^1 + 2$
 - $\rightarrow 2(x^3)^3 + 3(x^3)^2 + 1(x^3)^1 + 2$
 - $\rightarrow 2x^9 + 3x^6 + 1x^3 + 2$
 - $\rightarrow 2x^9 + 3x^6 + 1x^3 + 2 \ (x^4 = -1)$
 - $\rightarrow 2x + 17x^2 + 1x^3 + 2 \ (Z_{17})$
- Then perform key-switching from $s(X^t)$ to original s
- Require $2 * \log_2 n / 2$ (either direction) galois keys in Microsoft SEAL

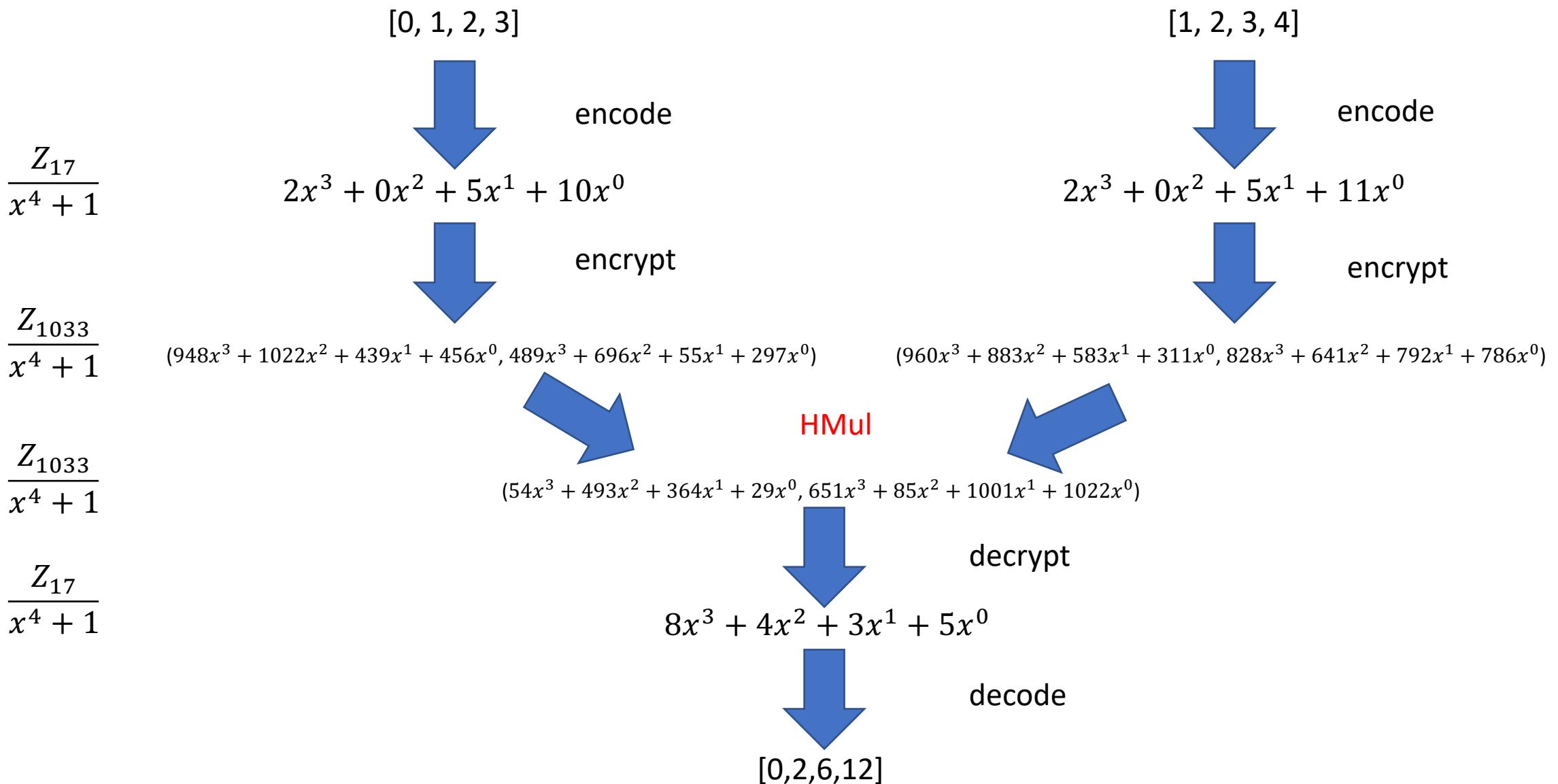
Routines



Routines



Routines



Routines

[157, 1014, 88, 516] [105, 86, 707, 13]

[0, 1, 2, 3]

$$\frac{Z_{17}}{x^4 + 1} \quad 2x^3 + 0x^2 + 5x^1 + 10x^0$$

encode



encrypt

$$\frac{Z_{17}}{x^4 + 1}$$

$$\frac{Z_{1033}}{x^4 + 1}$$

$$(948x^3 + 1022x^2 + 439x^1 + 456x^0, 489x^3 + 696x^2 + 55x^1 + 297x^0)$$

$$(516x^3 + 88x^2 + 1014x^1 + 157x^0, 13x^3 + 707x^2 + 86x^1 + 105x^0)$$

Rotate

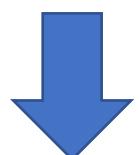


decrypt

$$\frac{Z_{1033}}{x^4 + 1}$$

$$\frac{Z_{17}}{x^4 + 1}$$

$$5x^3 + 0x^2 + 2x^1 + 10x^0$$

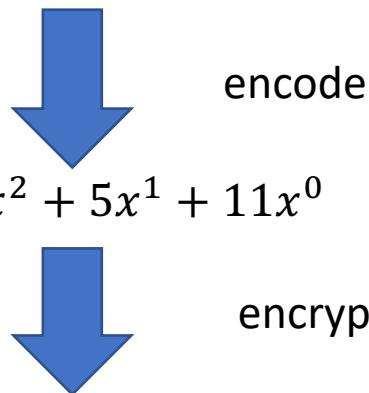


decode

[1, 0, 3, 2]

[1, 2, 3, 4]

$$2x^3 + 0x^2 + 5x^1 + 11x^0$$



encode

encrypt

$$(960x^3 + 883x^2 + 583x^1 + 311x^0, 828x^3 + 641x^2 + 792x^1 + 786x^0)$$

Rotate



decrypt



encode

$$(960x^3 + 883x^2 + 583x^1 + 311x^0, 828x^3 + 641x^2 + 792x^1 + 786x^0)$$

Summary

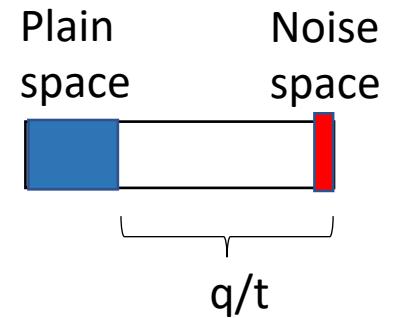
- Three main schemes
 - BGV
 - BFV (integer arithmetic)
 - CKKS (suitable for approximation arithmetic)
- Function call
 - **Encode, Decode**
 - **Encrypt, Decrypt**
 - Input a, b should be encrypted:
 - **H_Add, H_Mul, H_Rotate**
 - One input a, b can be plaintext
 - **PlainMul, PlainAdd**

Parameter selections

- t
 - Select first so that all plaintext operation is in t
 - Example: $t = 20\text{bits}$, $a=5\text{bits}$, can perform up to a^4

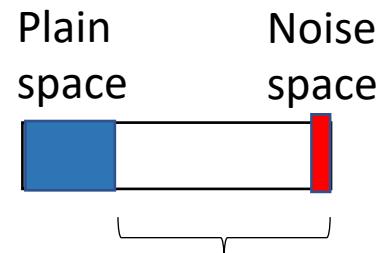
Parameter selections

- t
 - Select first so that all plaintext operation is in t
 - Example: $t = 20$ bits, $a=5$ bits, can perform up to a^4
- q
 - q/t should be large enough to ensure success of decryption



Parameter selections

- t
 - Select first so that all plaintext operation is in t
 - Example: t = 20bits, a=5bits, can perform up to a^4
- q
 - q/t should be large enough to ensure success of decryption
- n
 - $n/\log(q)$ should be large enough to ensure enough security level λ

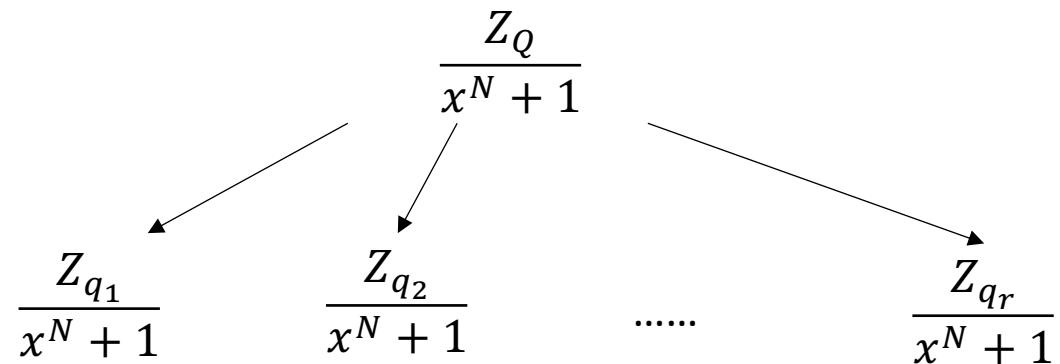


$$N \geq \frac{\lambda+110}{7.2} \log(Q) \text{ to get } \lambda\text{-bit security level}$$

n	Bit-length of default q		
	128-bit security	192-bit security	256-bit security
1024	27	19	14
2048	54	37	29
4096	109	75	58
8192	218	152	118
16384	438	300	237
32768	881	600	476

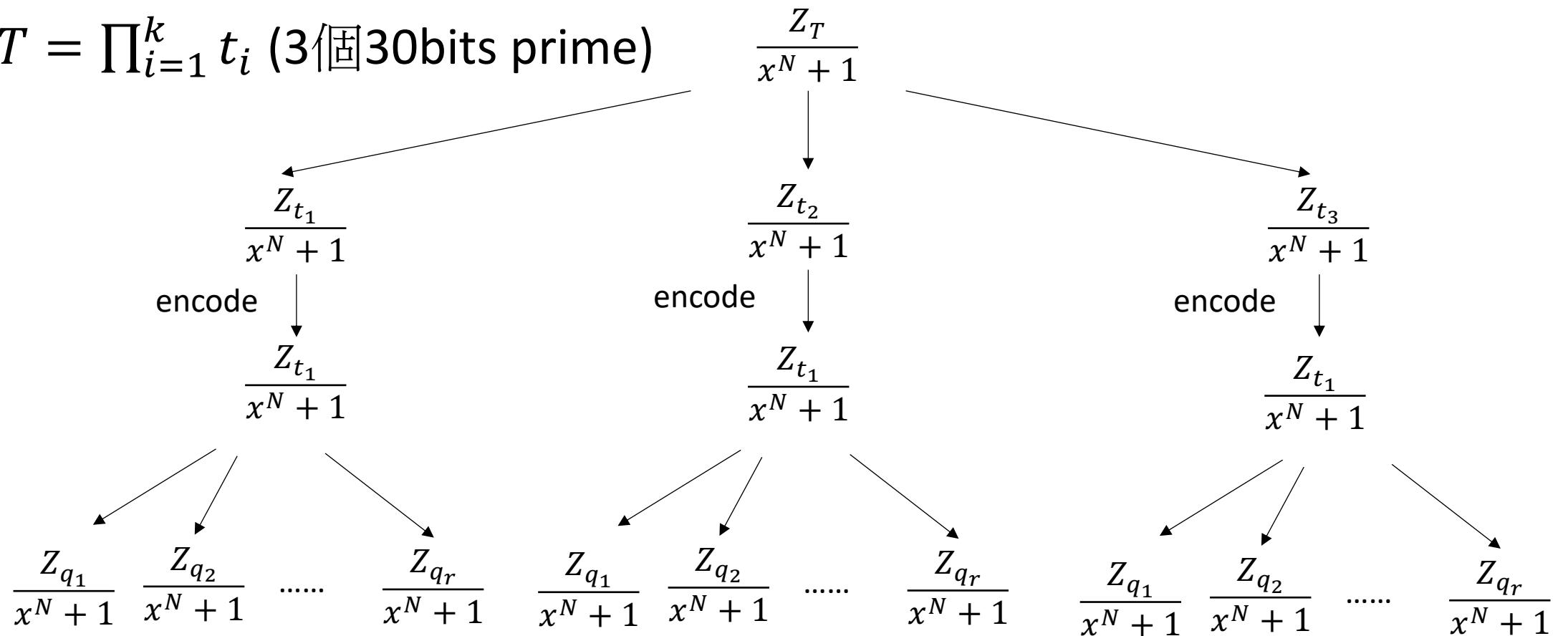
RNS (CRT)

- q may be huge (i.e. 300bits), which cannot fit in a machine word
- $Q = \prod_{i=1}^r q_i$ (10個30bits prime)
 - Operation in Q , equals operation in each q_i
 - More efficient, parallelizable



Double CRT and RNS

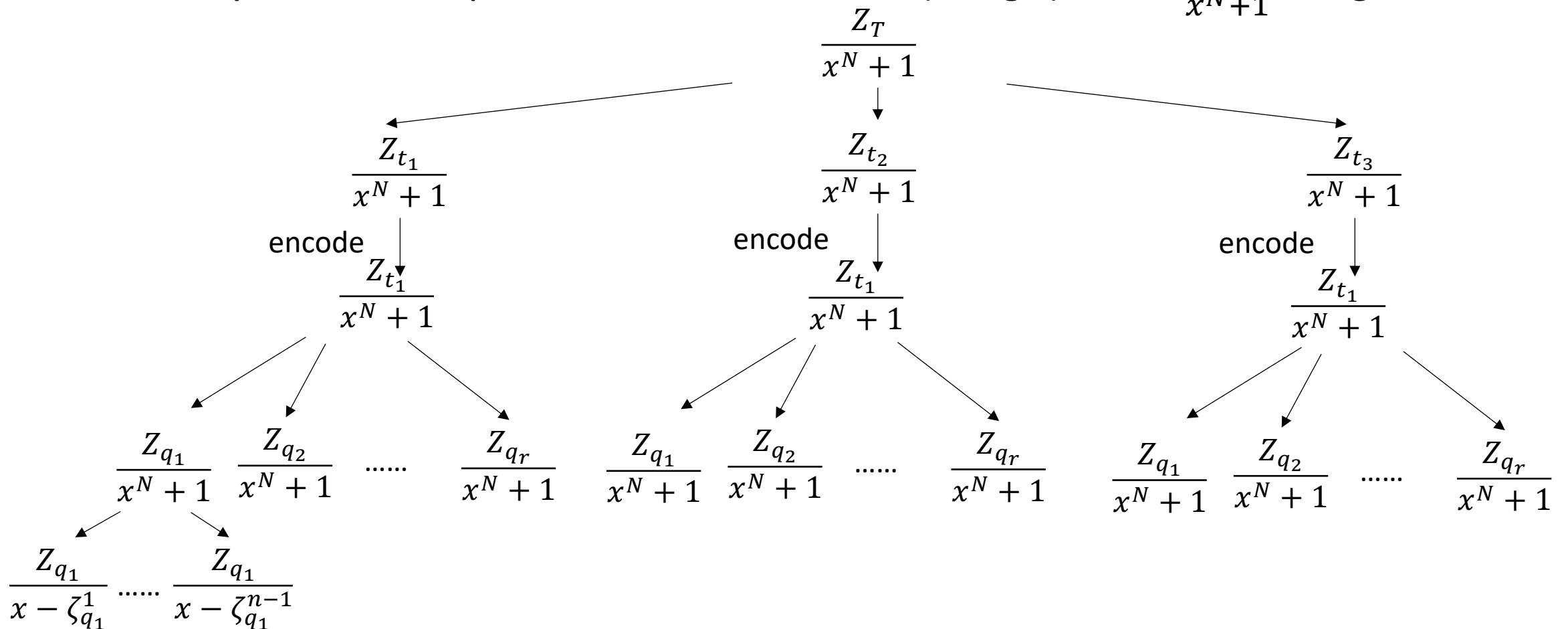
- t may also be huge (i.e. 90 bits)
- $T = \prod_{i=1}^k t_i$ (3個30bits prime)



Even further beyond

- Faster polynomial multiplication

- Polynomial multiplication can be done in $O(N \log N)$ under $\frac{Z_{q_1}}{x^N + 1}$ using NTT/FFT



Time, space, noise and security

Table 1. The asymptotic costs of homomorphic operations for the CKKS and RNS-CKKS scheme variants in HEAAN v1.0 and SEAL v3.1 respectively. $M(Q)$ is the complexity of multiplying large integers and is $O(\log^{1.58} Q)$ for HEAAN.

Homomorphic Operation	CKKS	RNS-CKKS with $Q = \prod_{i=1}^r Q_i$
addition, subtraction	$O(N \cdot \log Q)$	$O(N \cdot r)$
scalar multiplication	$O(N \cdot M(Q))$	$O(N \cdot r)$
plaintext multiplication	$O(N \cdot \log N \cdot M(O))$	$O(N \cdot r)$
ciphertext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$
ciphertext rotation	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$

Time, space, noise and security

Table 1. The asymptotic costs of homomorphic operations for the CKKS and RNS-CKKS scheme variants in HEAAN v1.0 and SEAL v3.1 respectively. $M(Q)$ is the complexity of multiplying large integers and is $O(\log^{1.58} Q)$ for HEAAN.

Homomorphic Operation	CKKS	RNS-CKKS with $Q = \prod_{i=1}^r Q_i$
addition, subtraction	$O(N \cdot \log Q)$	$O(N \cdot r)$
scalar multiplication	$O(N \cdot M(Q))$	$O(N \cdot r)$
plaintext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot r)$
ciphertext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$
ciphertext rotation	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$

Memory usage (bits) (RNS)

Private key:	2^*N
Public key:	$2^*\log_2 Q^*N$
Relinearization key:	$r * \log_2 Q^*N$
Rotation key:	$r * \log_2 Q^*N^*(\#rot)$

Time, space, noise and security

Table 1. The asymptotic costs of homomorphic operations for the CKKS and RNS-CKKS scheme variants in HEAAN v1.0 and SEAL v3.1 respectively. $M(Q)$ is the complexity of multiplying large integers and is $O(\log^{1.58} Q)$ for HEAAN.

Homomorphic Operation	CKKS	RNS-CKKS with $Q = \prod_{i=1}^r Q_i$
addition, subtraction	$O(N \cdot \log Q)$	$O(N \cdot r)$
scalar multiplication	$O(N \cdot M(Q))$	$O(N \cdot r)$
plaintext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot r)$
ciphertext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$
ciphertext rotation	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$

Operation	Time Cost	Noise Cost
Constant Addition	cheap	cheap
Addition	cheap	cheap
Constant Mult.	cheap	moderate ²
Multiplication	expensive	expensive
Rotation	expensive	cheap

Memory usage (bits) (RNS)

Private key:	2^*N
Public key:	$2^*\log_2 Q^*N$
Relinearization key:	$r * \log_2 Q^*N$
Rotation key:	$r * \log_2 Q^*N^*(\#rot)$

Time, space, noise and security

Table 1. The asymptotic costs of homomorphic operations for the CKKS and RNS-CKKS scheme variants in HEAAN v1.0 and SEAL v3.1 respectively. $M(Q)$ is the complexity of multiplying large integers and is $O(\log^{1.58} Q)$ for HEAAN.

Homomorphic Operation	CKKS	RNS-CKKS with $Q = \prod_{i=1}^r Q_i$
addition, subtraction	$O(N \cdot \log Q)$	$O(N \cdot r)$
scalar multiplication	$O(N \cdot M(Q))$	$O(N \cdot r)$
plaintext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot r)$
ciphertext multiplication	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$
ciphertext rotation	$O(N \cdot \log N \cdot M(Q))$	$O(N \cdot \log N \cdot r^2)$

Memory usage (bits) (RNS)	
Private key:	$2 * N$
Public key:	$2 * \log_2 Q * N$
Relinearization key:	$r * \log_2 Q * N$
Rotation key:	$r * \log_2 Q * N * (\#rot)$

Operation	Time Cost	Noise Cost
Constant Addition	cheap	cheap
Addition	cheap	cheap
Constant Mult.	cheap	moderate ²
Multiplication	expensive	expensive
Rotation	expensive	cheap

n	Bit-length of default q		
	128-bit security	192-bit security	256-bit security
1024	27	19	14
2048	54	37	29
4096	109	75	58
8192	218	152	118
16384	438	300	237
32768	881	600	476

Outline

- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

Comparison of existing encrypted inference

Table 2: MNIST performance comparison. Solutions are grouped by accuracy levels.

Method	Accuracy	Latency	
FHE-DiNN100	96.35%	1.65	(Bourse et al., 2017)
LoLa-Small	96.92%	0.29	
CryptoNets	98.95%	205	(Dowlin et al., 2016)
nGraph-HE	98.95% ⁹	135	(Boemer et al., 2018)
Faster-CryptoNets	98.7%	39.1	(Chou et al., 2018)
CryptoNets 2.3	98.95	24.8	
HCNN	99%	14.1	(Badawi et al., 2018)
LoLa-Dense	98.95%	7.2	
LoLa	98.95%	2.2	

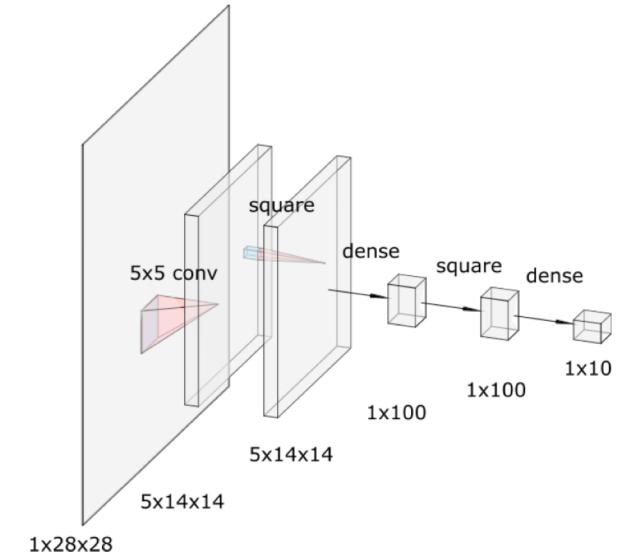


Figure 3: The structure of the network used for MNIST classification.

85625個double

Supported operations

$$\text{Ciphertext} \times \text{Scalar} = \text{Ciphertext} \quad \boxed{2|3|7|\dots} \times \boxed{2|2|2|\dots} = \boxed{4|6|14|\dots}$$

\downarrow
 ct

\downarrow
 plaintext scalar

\downarrow
 ct

$$\text{Ciphertext} \times \text{Plaintext} = \text{Ciphertext} \quad \boxed{2|3|7|\dots} \times \boxed{1|2|3|\dots} = \boxed{2|6|21|\dots}$$

\downarrow
 ct

\downarrow
 pt

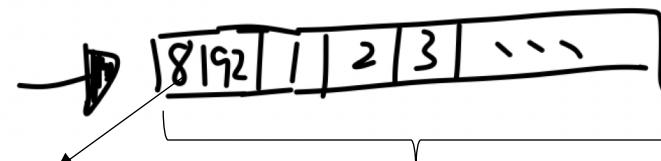
$$\text{Ciphertext} \times \text{Ciphertext} = \text{Ciphertext} \quad \boxed{2|3|7|\dots} \times \boxed{1|2|3|\dots} = \boxed{2|6|21|\dots}$$

\downarrow
 ct

\downarrow
 pt

加法一樣 (elementwise)

Rotate 



80-bit integer

8192 slots

Scenario

- Inference MNIST (28x28)

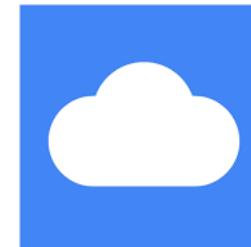
Here,
 $N=8192$, $k=4$, $\lambda = 128$
 $Q=\{43,43,44,44,44\}$ bits, $r=5$
 $\log_2 Q \approx 218$
 $\#rot \approx 20$

Memory usage (bits) (RNS)

Private key:	$2*N$	2KB
Public key:	$2*\log_2 Q*N$	420KB
Relinearization key:	$r * \log_2 Q*N$	1MB
Rotation key:	$r * \log_2 Q*N*(\#rot)$	20MB
Encrypted 28x28 image:	$k*25*2*\log_2 Q*N$	42MB
Result returned:	$2*\log_2 Q*N$	420KB

$420KB+1MB+20MB+42MB$

$\approx 64MB$



420KB

Scenario

- Inference MNIST (28x28)

Here,
 $N=8192$, $k=4$, $\lambda = 128$
 $Q=\{43,43,44,44,44\}$ bits, $r=5$
 $\log_2 Q \approx 218$
 $\#rot \approx 20$

Memory usage (bits) (RNS)

Private key:	$2*N$	2KB
Public key:	$2*\log_2 Q*N$	420KB
Relinearization key:	$r * \log_2 Q*N$	1MB
Rotation key:	$r * \log_2 Q*N*(\#rot)$	20MB
Encrypted 28x28 image:	$k*25*2*\log_2 Q*N$	42MB
Result returned:	$2*\log_2 Q*N$	420KB

$420KB+1MB+20MB+42MB$

$\approx 64MB$



420KB

Cf:
RSA-2048, public key 2048bits
Message 245B → 256B

Neural network

- Cannot perform ReLU or sigmoid straightforward
 - Approximation using taylor expansion
 - Using HE friendly activation function
 - Ex: x^2

Neural network

- Cannot perform ReLU or sigmoid straightforward
 - Approximation using taylor expansion
 - Using HE friendly activation function
 - Ex: x^2

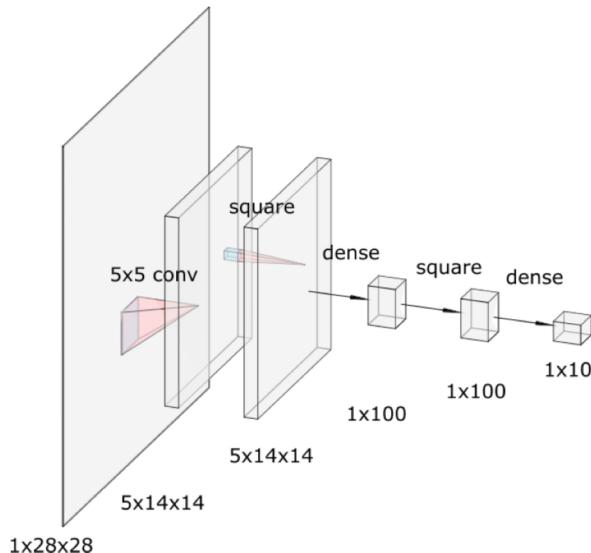


Figure 3: The structure of the network used for MNIST classification.

Neural network

- Cannot perform ReLU or sigmoid straightforward
 - Approximation using taylor expansion
 - Using HE friendly activation function
 - Ex: x^2

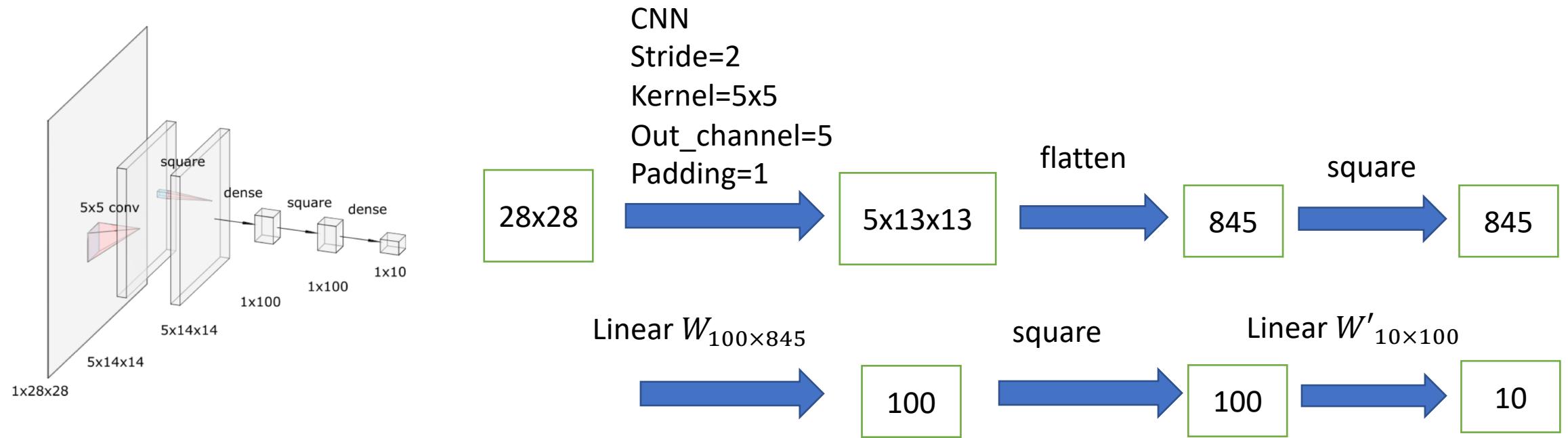


Figure 3: The structure of the network used for MNIST classification.

Training phase

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(1, 5, 5, 2, padding=1)
        self.fc1 = nn.Linear(5*169, 100)
        self.fc2 = nn.Linear(100, 10)

    def forward(self, x):
        x = self.conv1(x)
        x = x**2
        x = torch.flatten(x, 1)
        x = self.fc1(x)
        x = x**2
        x = self.fc2(x)
        output = F.log_softmax(x, dim=1)
        return output
```

Test set: Average loss: 0.0837, Accuracy: 9874/10000 (99%)

Inference

- Take the model parameters out
- Using fixed point to represent model coefficients and inputs
- Divide the result with scale afterwards

Here,
 $N=8192, k=4, \lambda = 128$
 $Q=\{43,43,44,44,44\}$ bits, $r=5$
 $\log_2 Q \approx 218$
 $\#rot \approx 20$

CNN

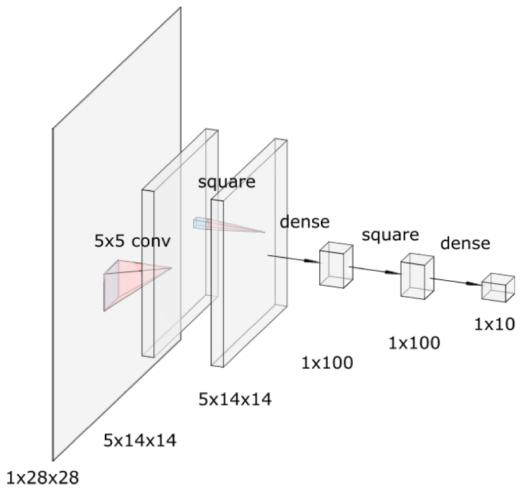
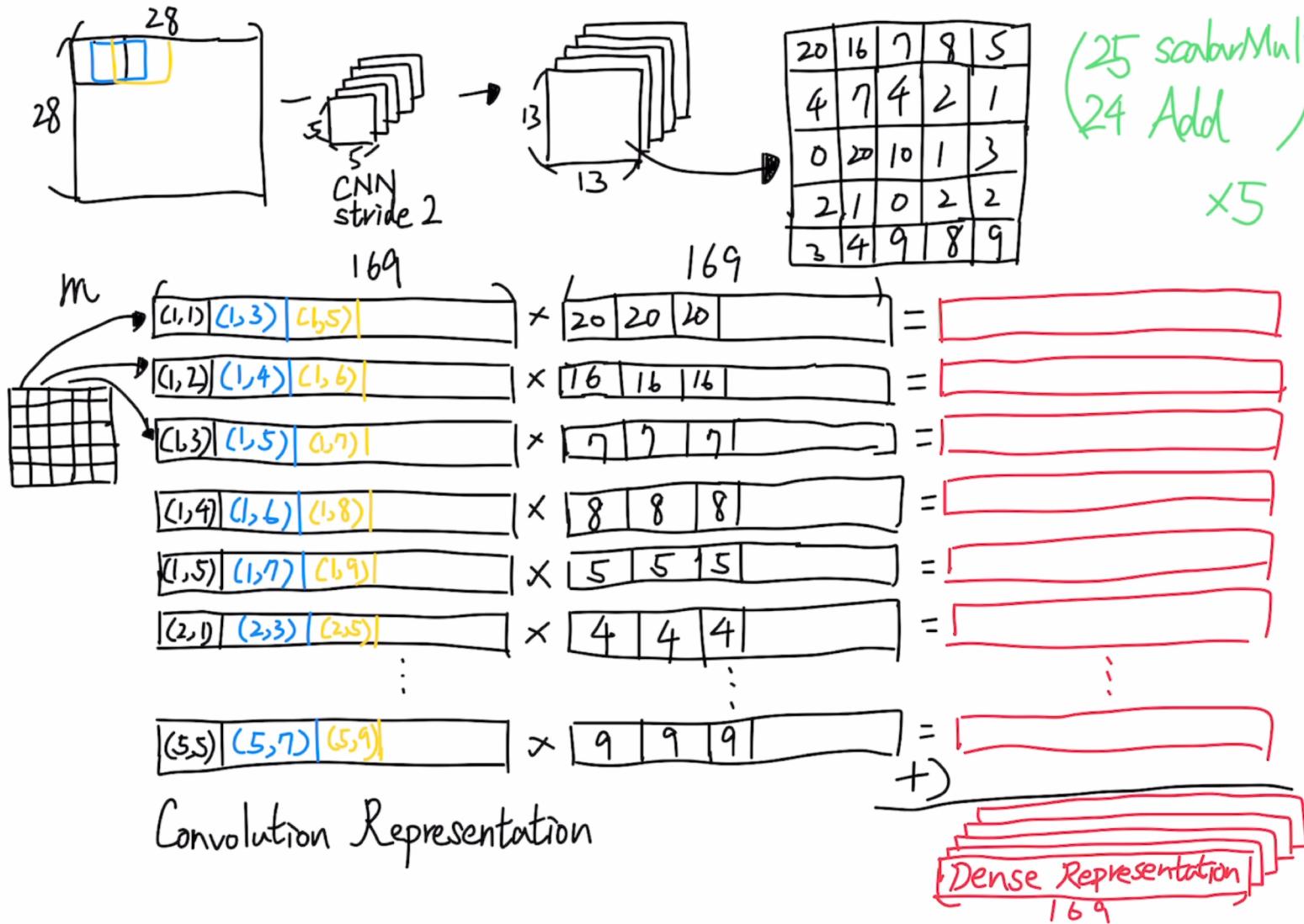
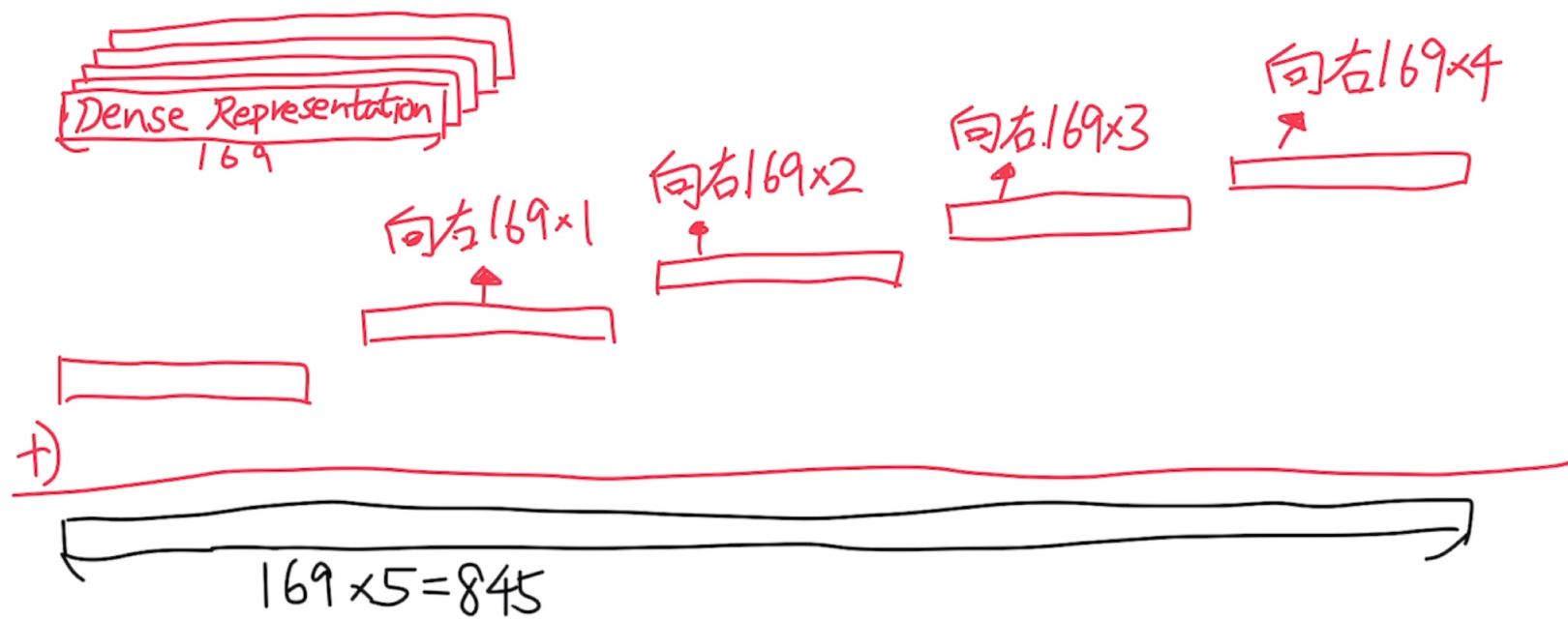
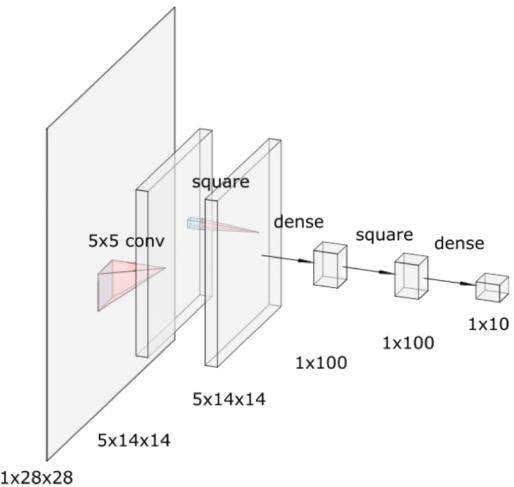


Figure 3: The structure of the network used for MNIST classification.

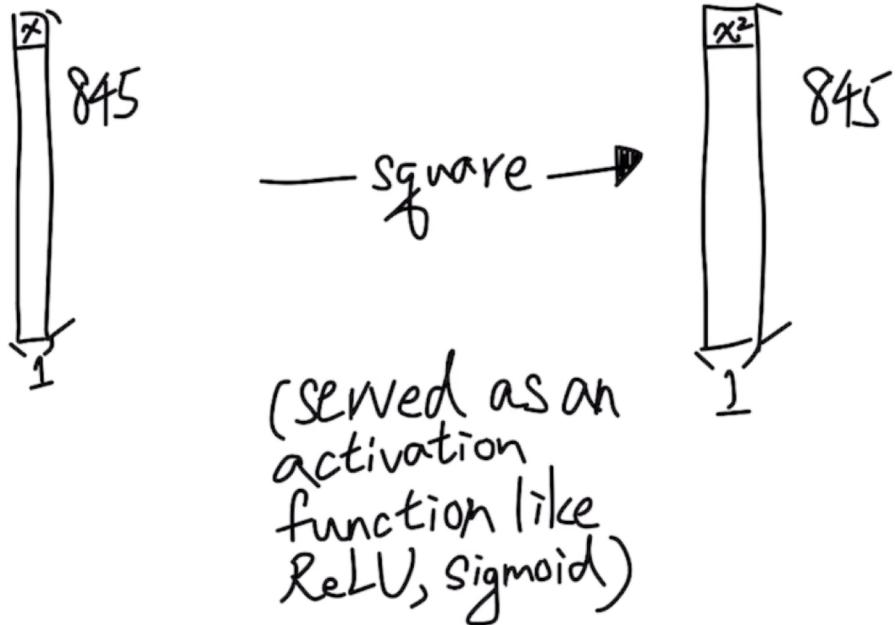
Flatten



4 Rotation
4 Add



Square



1 cipherMu

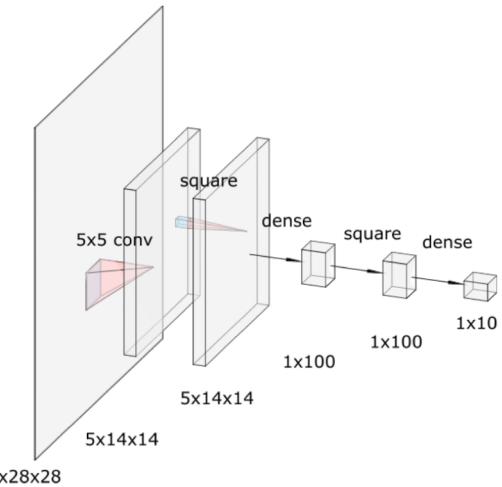


Figure 3: The structure of the network used for MNIST classification.

Dense layer 845 → 100

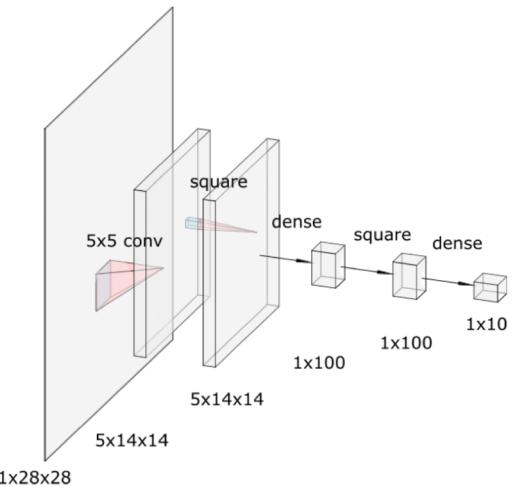
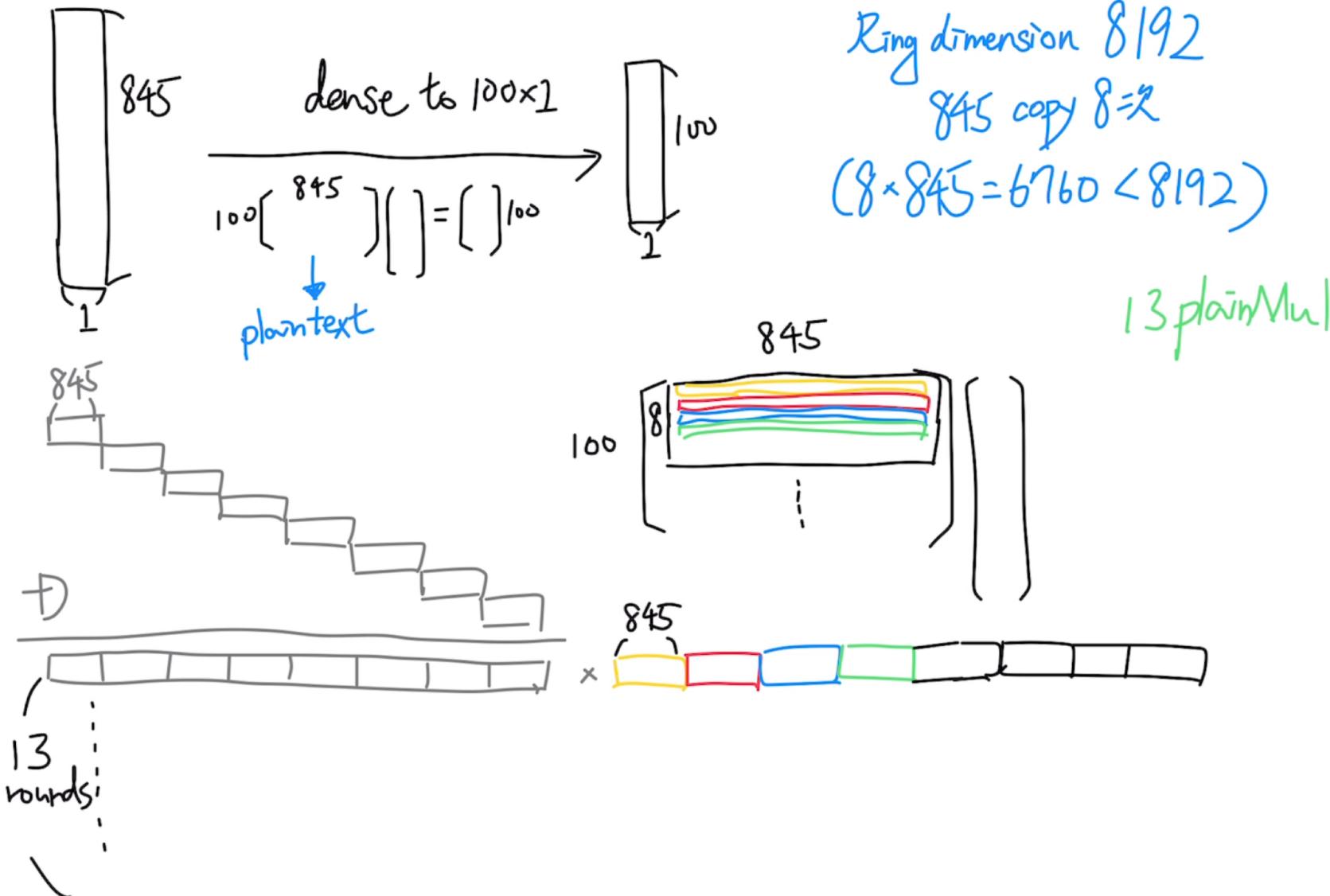
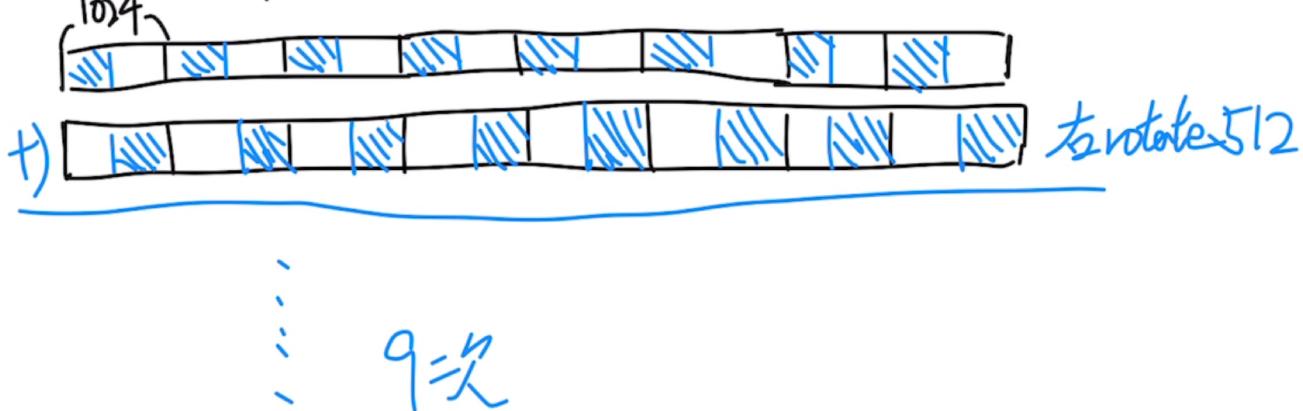


Figure 3: The structure of the network used for MNIST classification.

Dense layer 845 → 100

實際上
 845
 1024
 1024

Rotate 9
 Add 9



最後變成

最左是
reduce sum 結果

13
 這13個再
combine一起
 12Rotation
 12Add
 interleaving representation

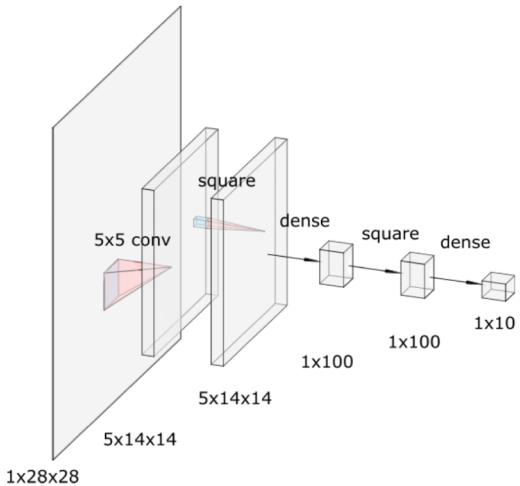
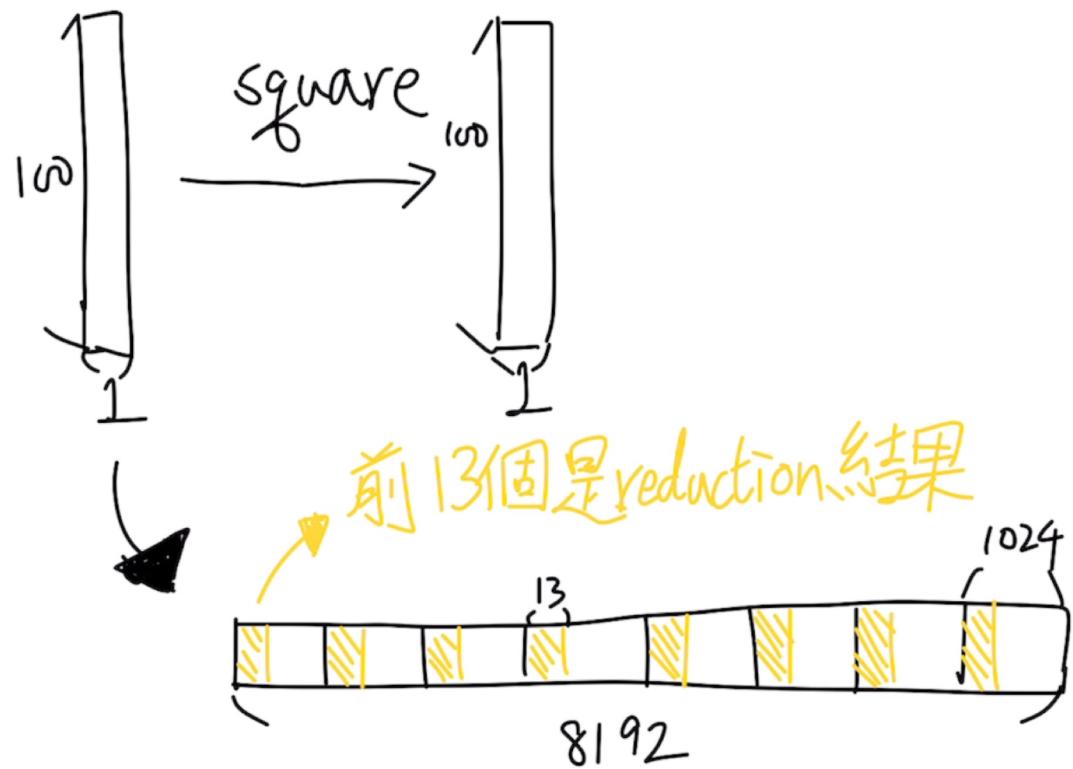


Figure 3: The structure of the network used for MNIST classification.

Square



$$10 \begin{bmatrix} & \\ & \end{bmatrix} = \begin{bmatrix} & \\ & 100 \end{bmatrix}$$

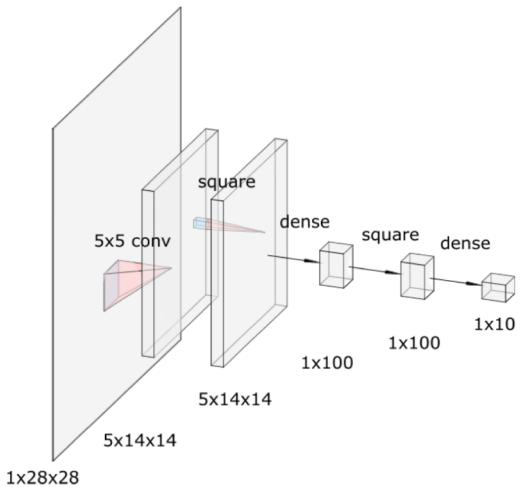


Figure 3: The structure of the network used for MNIST classification.

Dense layer 100 → 10

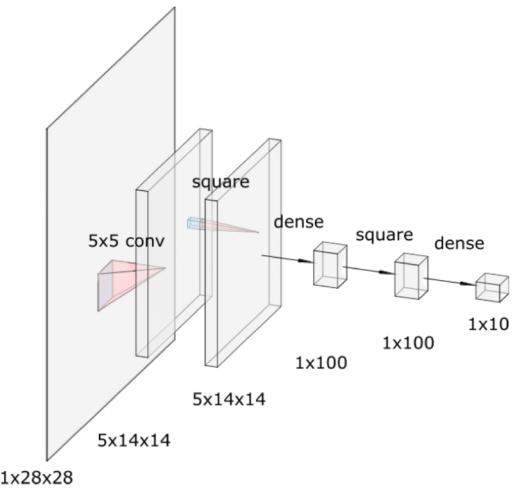
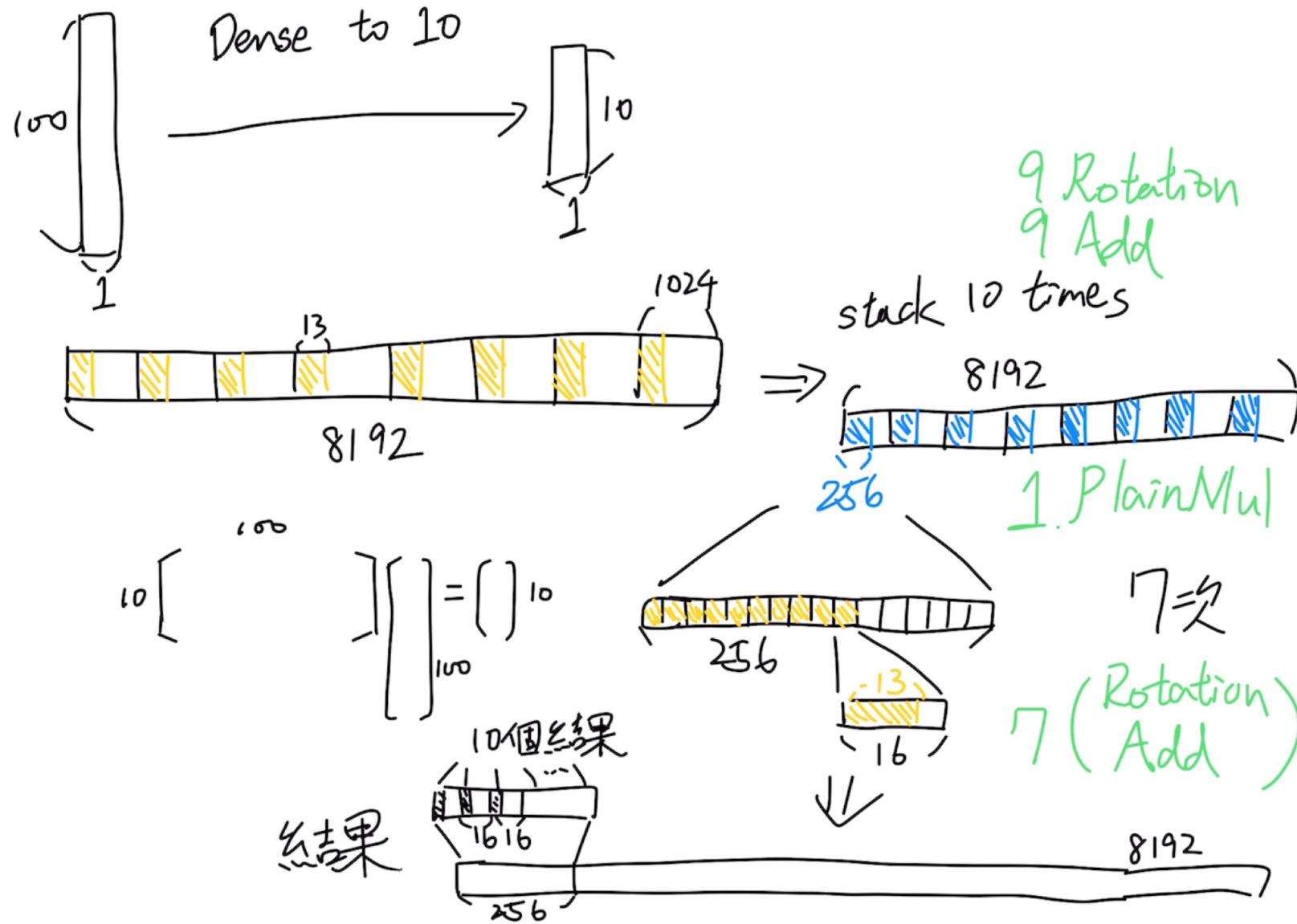
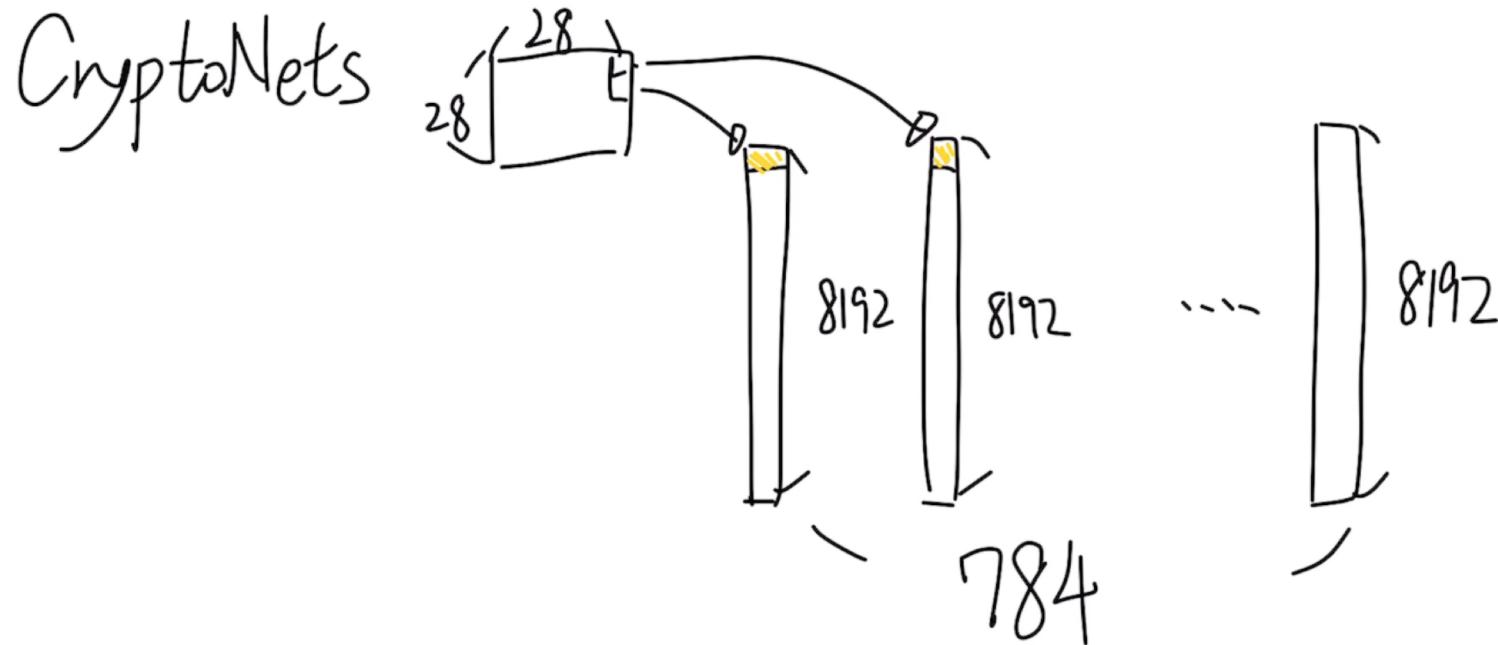


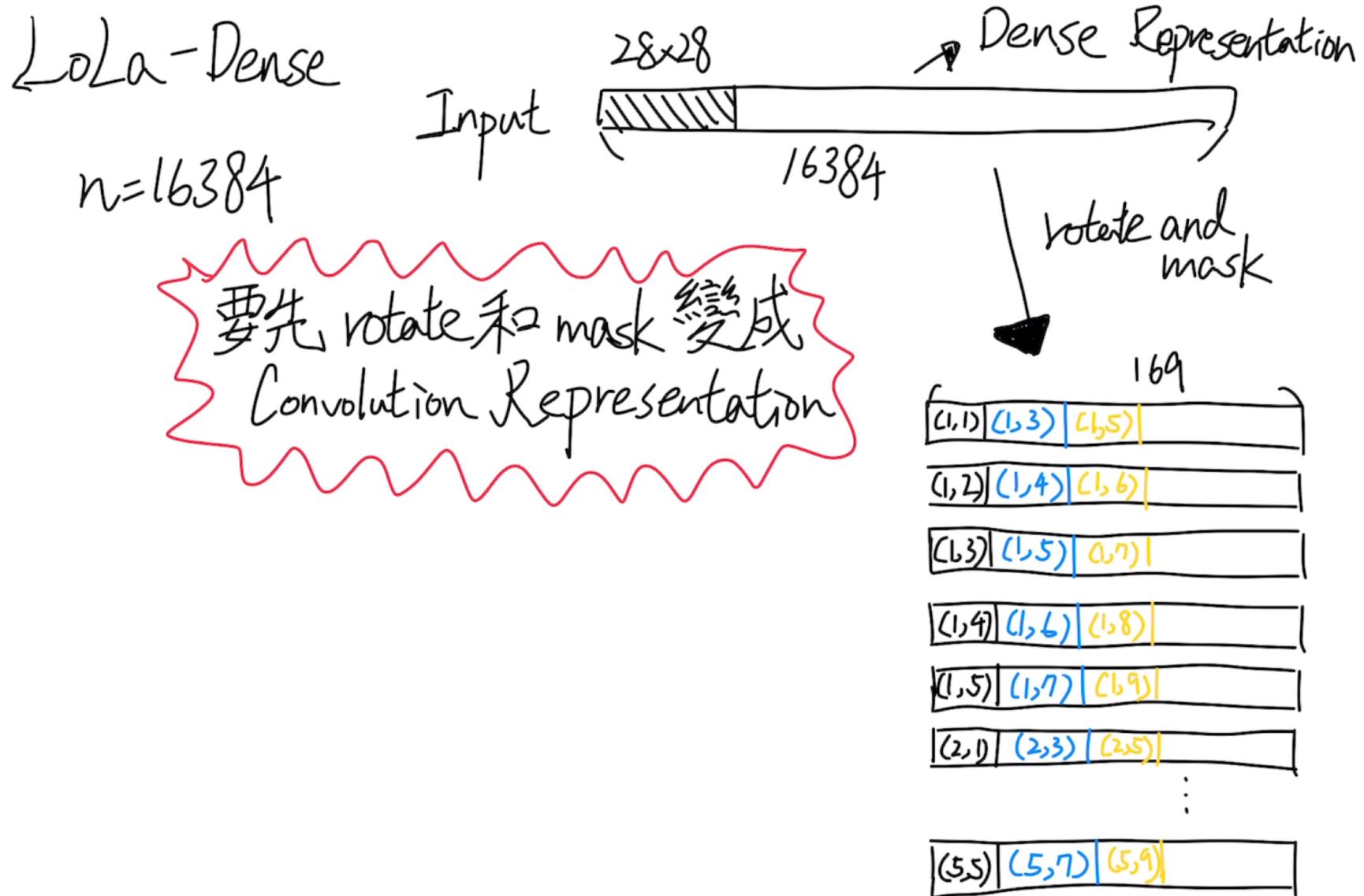
Figure 3: The structure of the network used for MNIST classification.

CryptoNets



High-latency, but can perform
batch-operation ($bs = 8192$)
(amortized latency)

LoLa-Dense



LoLa-Small

Lola small

One CNN → square activation
→ dense to 10

$$125 + 845 \times 10 = 8575 \text{ doubles}$$

Comparison

Azure standard B8ms virtual machine with 8 vCPUs and 32GB of RAM

Table 2: MNIST performance comparison. Solutions are grouped by accuracy levels.

Method	Accuracy	Latency	
FHE-DiNN100	96.35%	1.65	(Bourse et al., 2017)
LoLa-Small	96.92%	0.29	
CryptoNets	98.95%	205	(Dowlin et al., 2016)
nGraph-HE	98.95% ⁹	135	(Boemer et al., 2018)
Faster-CryptoNets	98.7%	39.1	(Chou et al., 2018)
CryptoNets 2.3	98.95	24.8	
HCNN	99%	14.1	(Badawi et al., 2018)
LoLa-Dense	98.95%	7.2	
LoLa	98.95%	2.2	

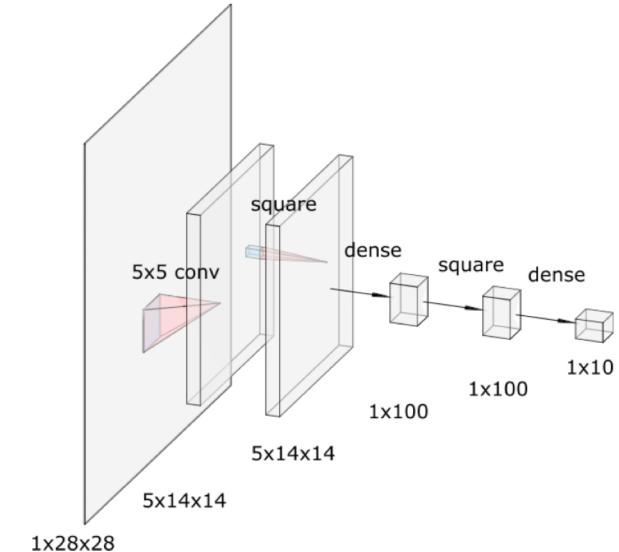


Figure 3: The structure of the network used for MNIST classification.

85625個double

11x faster than CryptoNets 2.3 and 93x faster than CryptoNets

Scaling

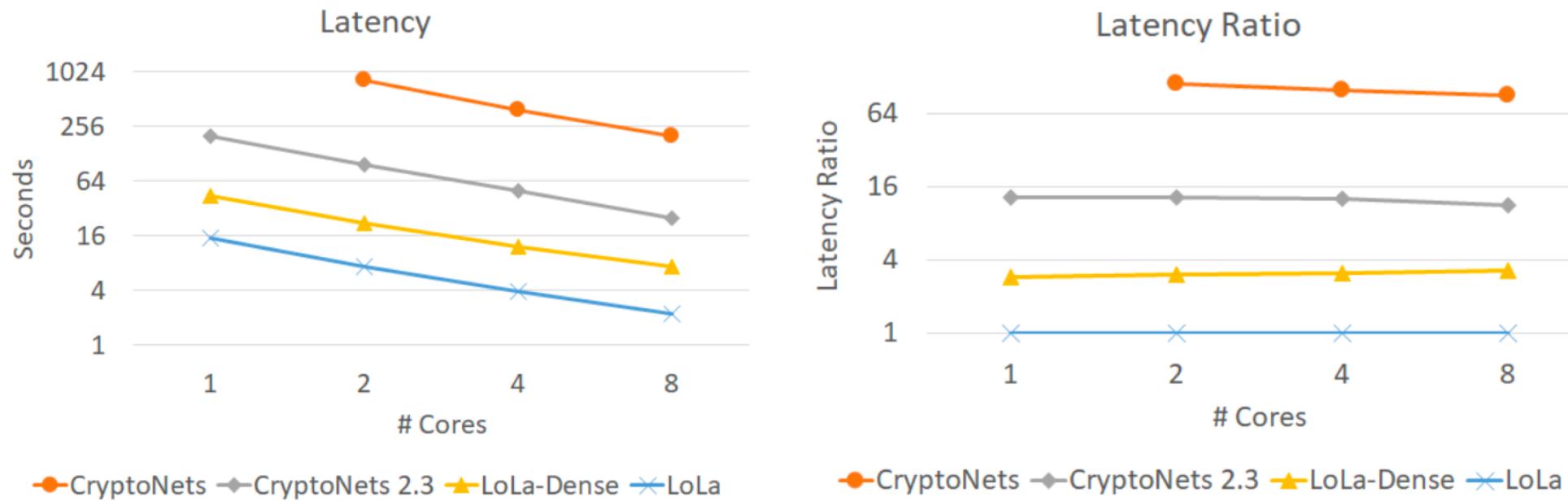


Figure 2: Left: latency of the different network implementations for the MNIST task with respect to the number of available cores. Right: ratio between the latency of each solution and the latency of LoLa

Testing

- Intel(R) Core(TM) i7-6700K CPU @ 4.00GHz, 8 cores, 64GB RAM

```
+-----+
|       BFV Performance Test with Degree: 8192      |
+-----+
/
| Encryption parameters :
|   scheme: BFV
|   poly_modulus_degree: 8192
|   coeff_modulus size: 218 (43 + 43 + 44 + 44 + 44) bits
|   plain_modulus: 786433
\
```

```
Average batch: 123 microseconds
Average unbatch: 152 microseconds
Average encrypt: 4859 microseconds
Average decrypt: 1395 microseconds
Average add: 60 microseconds
Average multiply: 14989 microseconds
Average multiply plain: 2456 microseconds
Average square: 10842 microseconds
Average relinearize: 5398 microseconds
Average rotate rows one step: 5425 microseconds
Average rotate rows random: 24049 microseconds
Average rotate columns: 5414 microseconds
```

μs	Without HE	With HE
CNN	222	MulPlain:125 Add:129 Rotate:4 $=336440$
Square_1	9	CipherMul: 1 10842
Dense to 100	952	Rotate: 136 Add: 149 MulPlain: 13 $=778668$
Square_2	9	CipherMul: 1 10842
Dense to 10	12	Rotate: 16 Add: 17 MulPlain: 16 $(9+7)*(5480)+2456$ $=90292$
Total	1204	$4*1685488$ $=4908336$

Outline

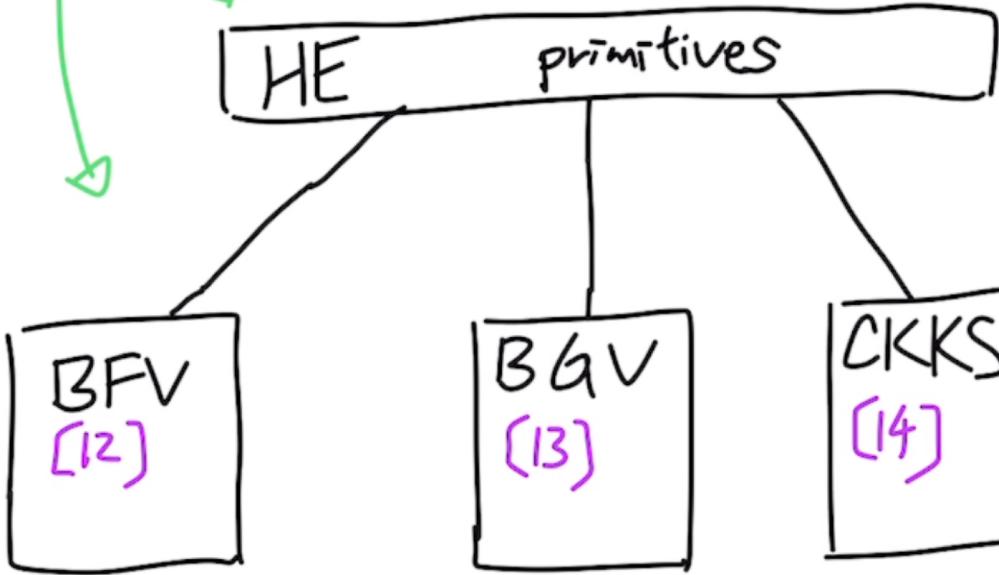
- Introduction
- Possible scenarios
- More specific
 - BFV encryption
 - Encrypted inference (LoLa)
- Research directions

FL +
Hybrid
scheme
[21], [22]

HE

Compiler

Applications, e.g.- GWAS [1], Private Query [2],
Inference [3], [4], [5], [6], [7],
Training [8]



Add
Mul
Rot
Enc
Dec

Scheme acceleration [15], [16], [17]

GPU [18], [19], [20]
FPGA [23]
ASIC [24]
RISC-HE

Reference

- On the explanation of FHE: https://www.youtube.com/watch?v=pXb39wj5ShI&ab_channel=OfficeoftheDirectorofNationalIntelligence
- Possible HE application: https://www.researchgate.net/publication/320976976_APPLICATIONS_OF_HOMOMORPHIC_ENCRYPTION
- Covid-19 contact tracing: https://dualitytech.com/wp-content/uploads/2020/04/Covid19_onepager.pdf
- Open source HE library: <https://scholarworks.uark.edu/cgi/viewcontent.cgi?article=1074&context=csceuht>
- Somewhat Homomorphic encryption: <https://eprint.iacr.org/2012/144.pdf>
- Homomorphic Encryption for Arithmetic of Approximate Numbers: <https://eprint.iacr.org/2016/421.pdf>
- Faster Homomorphic Linear Transform in HELib: <https://eprint.iacr.org/2018/244.pdf>
- Homomorphic Evaluation of the AES Circuit: <https://eprint.iacr.org/2012/099.pdf>
- SEAL-manual: <https://www.microsoft.com/en-us/research/uploads/prod/2017/11/sealmanual-2-3-1.pdf>
- CHET: <https://www.cs.utexas.edu/~roshan/CHET.pdf>
- Multi-GPU RNS: <https://ieeexplore.ieee.org/document/9185077>