



浙江大学 区块链与数据安全  
全国重点实验室  
STATE KEY LABORATORY OF BLOCKCHAIN AND DATA SECURITY  
ZHEJIANG UNIVERSITY

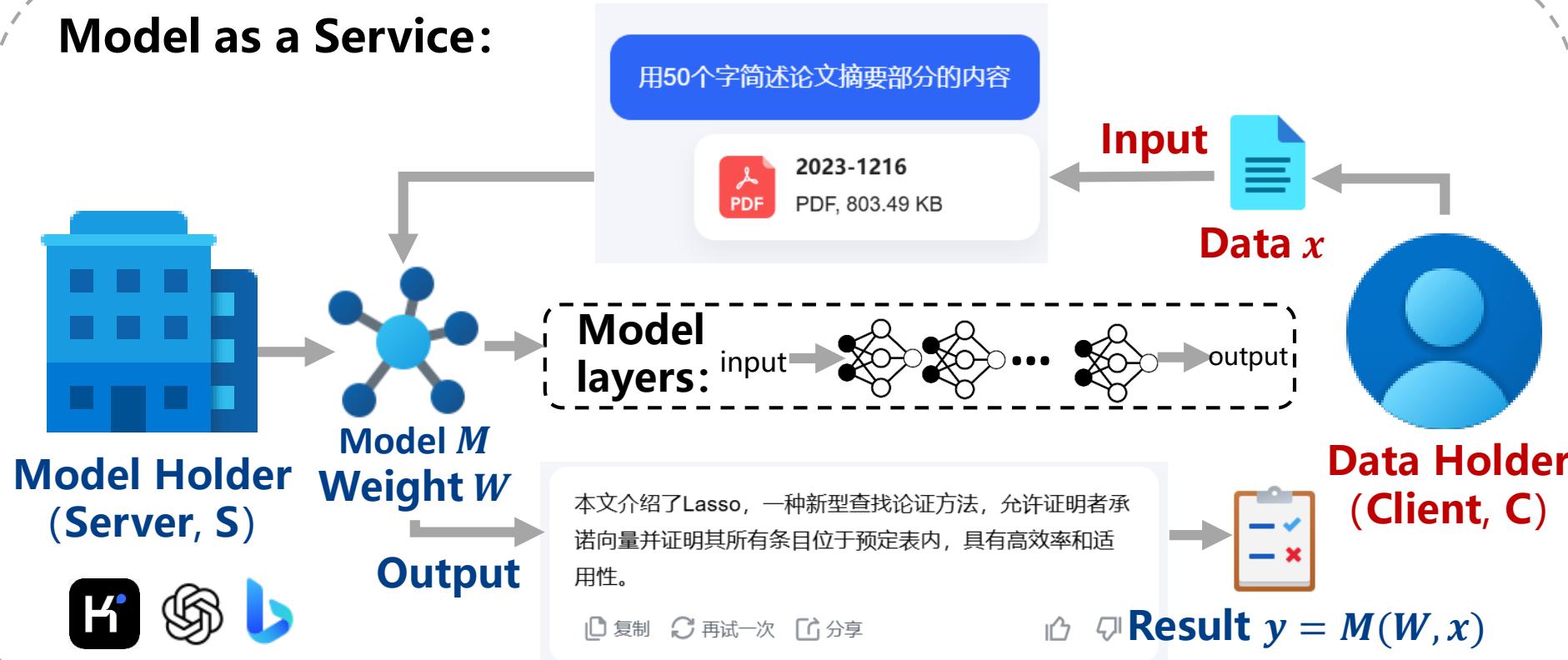
# A New PPML Paradigm for Quantized Models

Tianpei Lu, Bingsheng Zhang\*, Xiaoyuan Zhang, Kui Ren

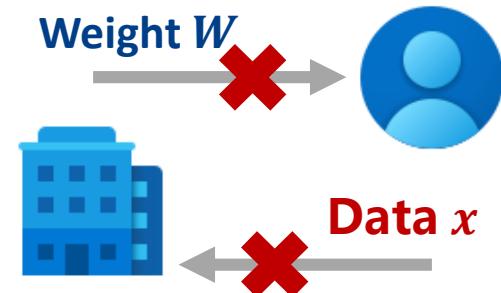
Zhejiang University

# Background

## Model as a Service:



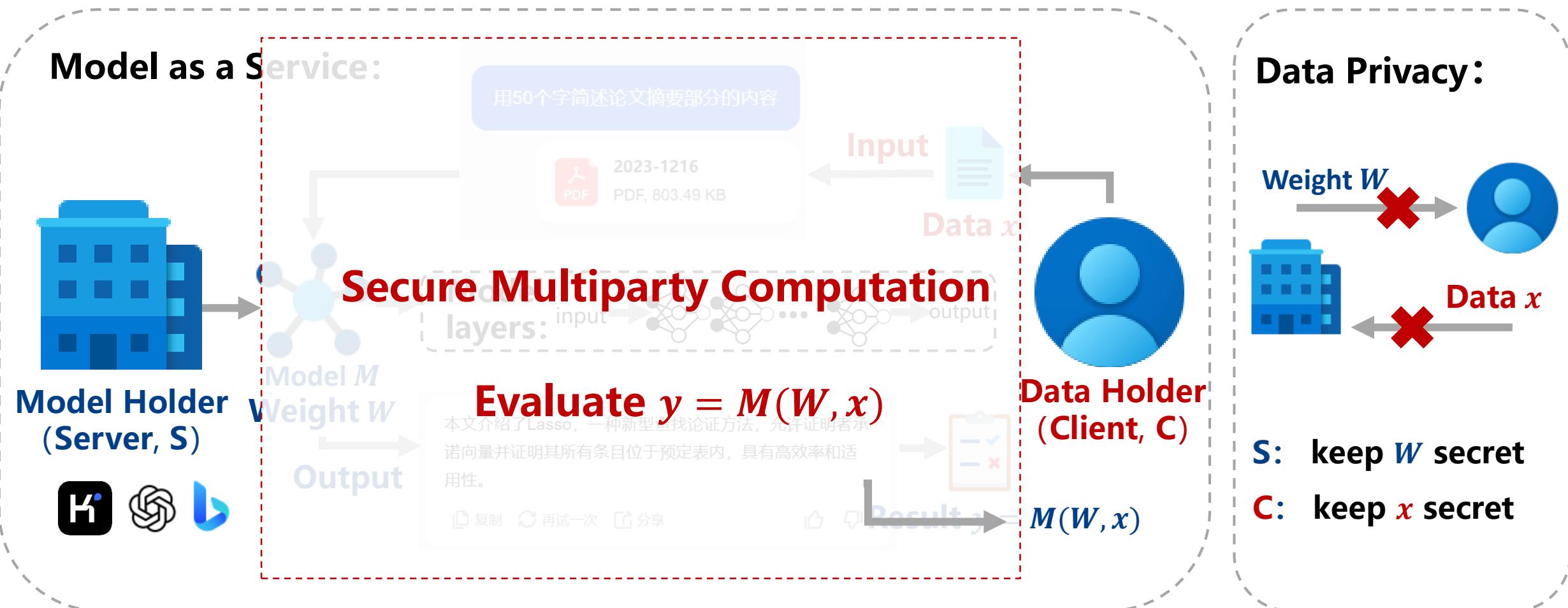
## Data Privacy:



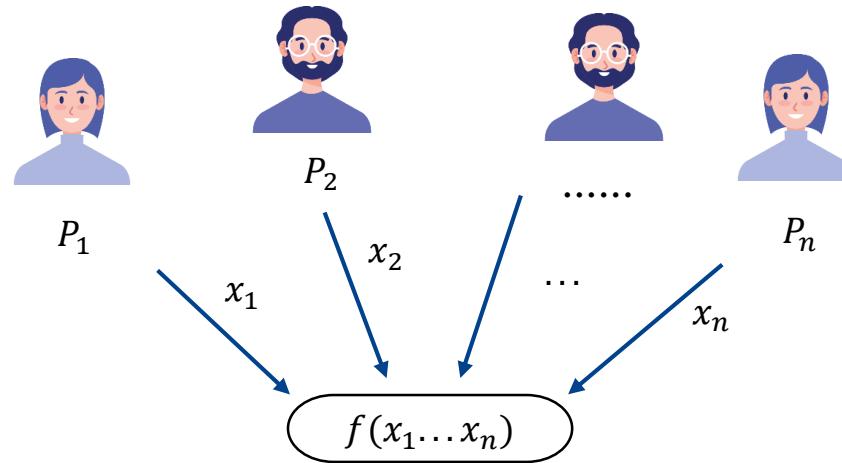
**S:** keep  $W$  secret

**C:** keep  $x$  secret

# Background



# Secure Multiparty Computation



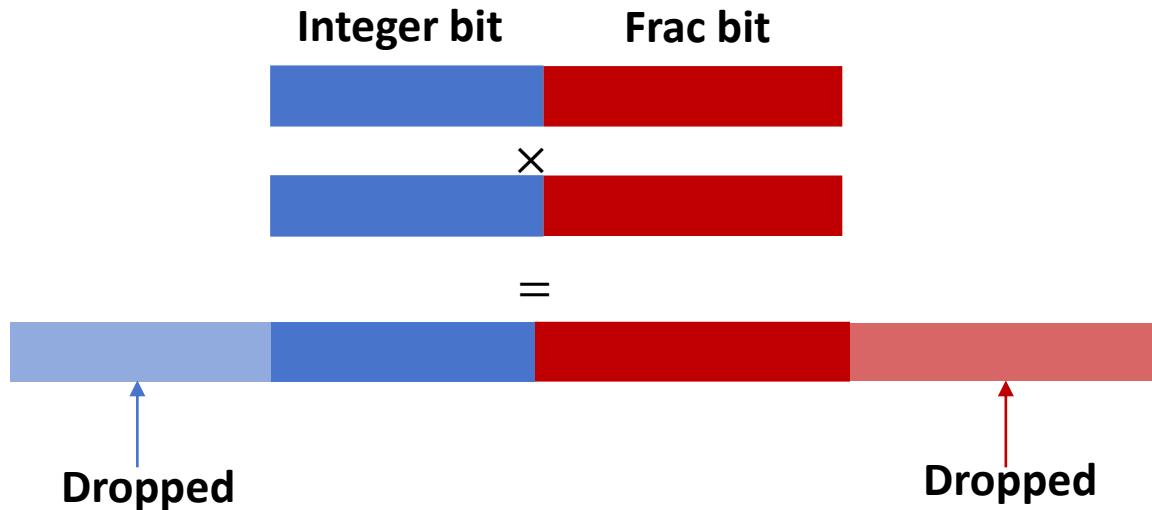
$P_i$  input messages  $x_i$  and receive  $y_i$

$$y_1 \dots y_n = f(x_1 \dots x_n)$$

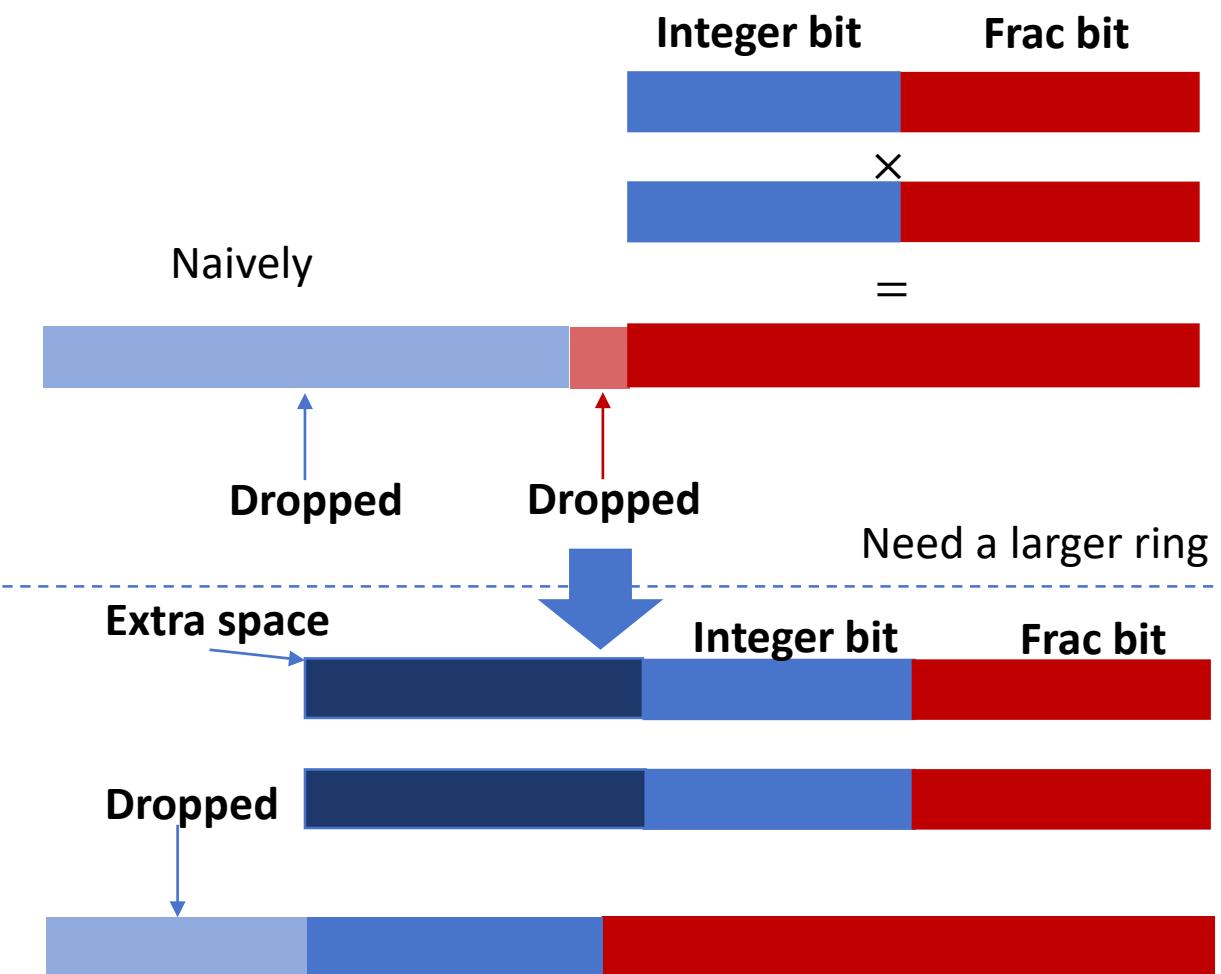
Typical Solution: Secure Multiparty Computation

# Challenges of PPML (1)

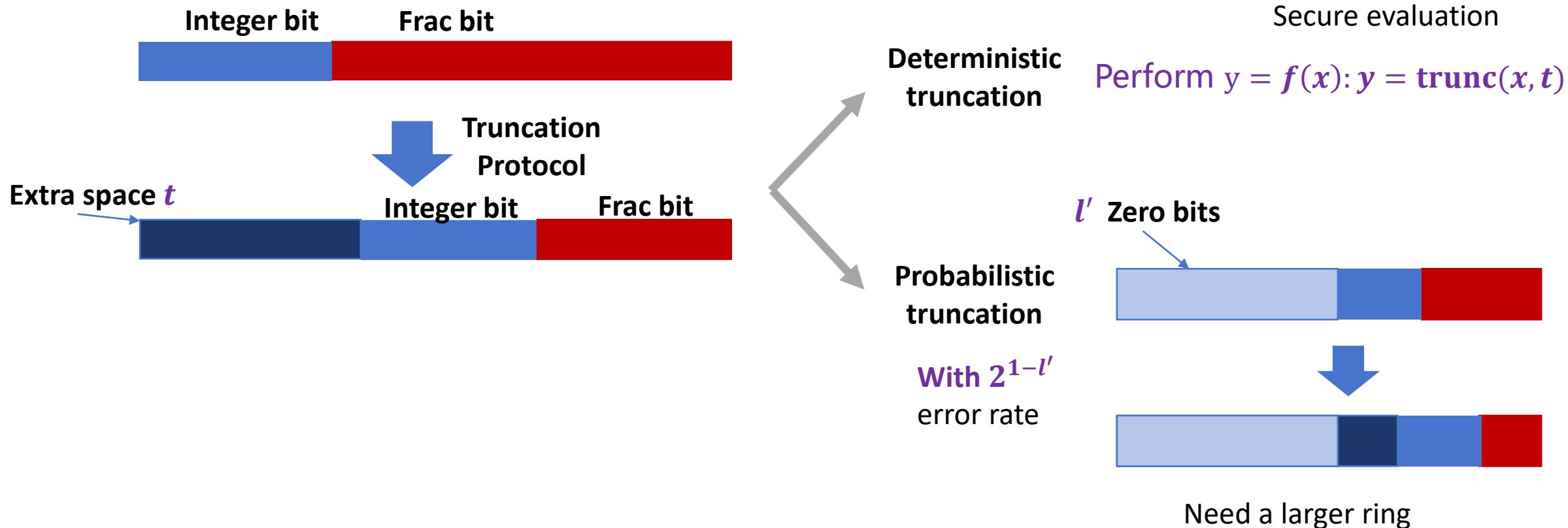
## Fixed point multiplication



## Multiplication over ring

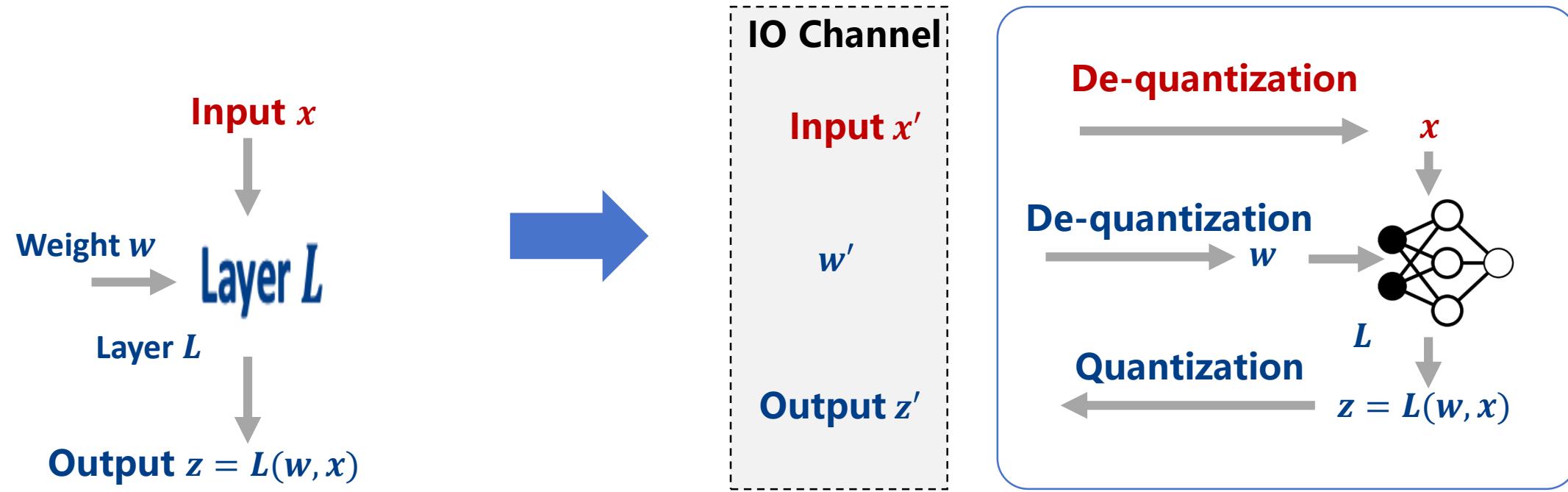


# Challenges of PPML (2)



Truncation is expensive !!!

# Quantization in PPML



$x$ 、 $w$ 、 $z$  are Fp32 or Fp64

$x'$ 、 $w'$ 、 $z'$  are Int8

$$\text{Quantization: } X = \frac{x'}{s_X} + b_X$$

$$\text{De-quantization: } X = s_X(X' - b_X)$$

# Quantization in PPML

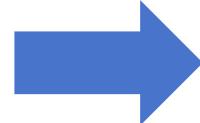
## Challenges in Quantization PPML

$$X' = \frac{X}{s_X} + b_X$$

$$Y' = \frac{Y}{s_Y} + b_Y$$

$$Z = XY$$

$$Z' = \frac{Z}{s_Z} + b_Z$$



$$Z' = \frac{(X' - b_X)(Y' - b_Y)}{\frac{s_X s_Y}{s_Z}} + b_Z$$

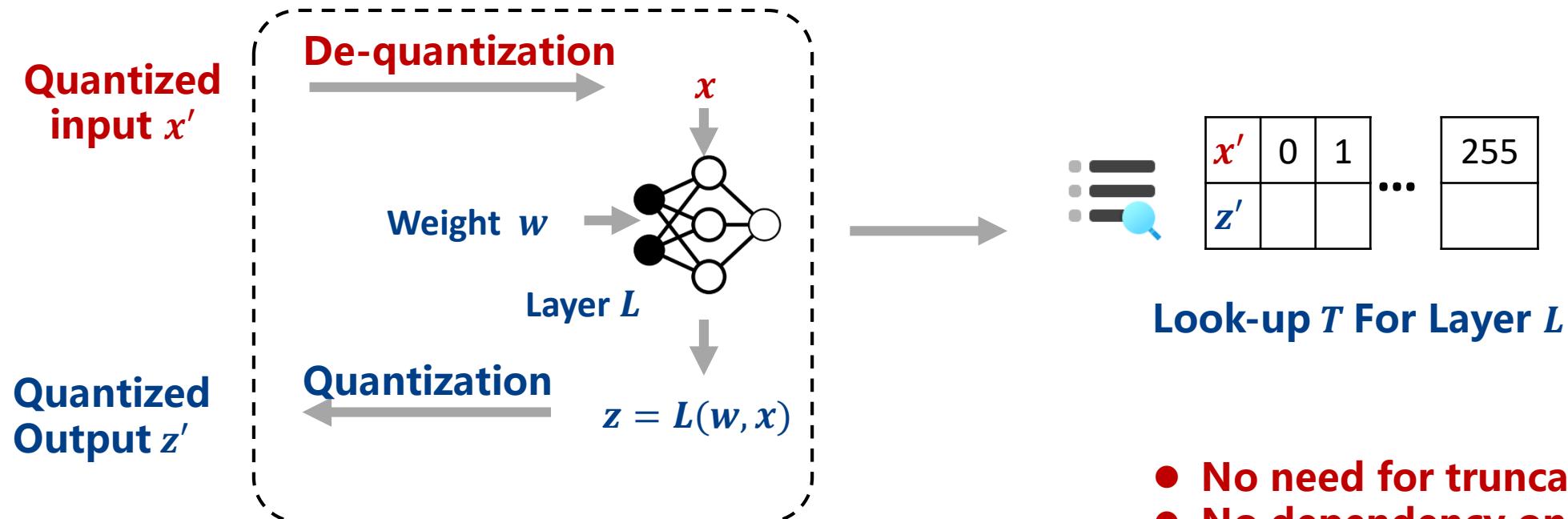
Int8 multiplication:  
Calculated over Int16

Fixed point scaler:  
Require truncation

Int8 offset

# Quantization & Look-up table

Look-up table evaluation is all we need:



- No need for truncation
- No dependency on internal computational complexity.

# Private look-up evaluation protocol

$$x = [x]_0 + [x]_1$$

Evaluate  $[x] \rightarrow [f(x)]$  using Look-up table

Input  $[x] := \{[x]_0, [x]_1\}$

Shifted look-up table triple  $\{[r], [f(r)], \dots, [f(r + n - 1)]\}$

0	$[f(r)]_0$
1	$[f(r + 1)]_0$
.....	
$x - r$	$[f(x)]_0$
.....	
$n - 1$	$[f(r + n - 1)]_0$



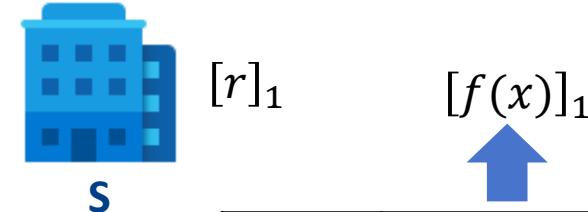
$[r]_0$

$[f(x)]_0$

1. Open  $x - r$   
 2. select  $x - r$  th item

$O(1)$  online communication

**Server knows  $f$ , but should not know  $r$ !!!**



0	$[f(r)]_1$
1	$[f(r + 1)]_1$
.....	
$x - r$	$[f(x)]_1$
.....	
$n - 1$	$[f(r + n - 1)]_1$

# Private look-up evaluation protocol

$(n-1, n)$ -ROT :

$n$   $l$ -bit messages  
 $(m_0, \dots, m_n)$



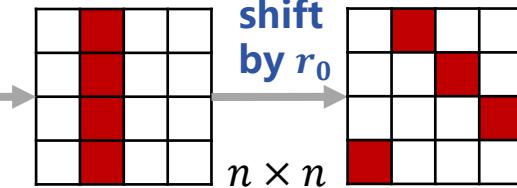
Sender

E.g.: can be constructed from  $\log n$  ROT via GGM tree

Receiver

$$M = (m_0, \dots, m_j, \dots, m_N)$$

Left shift by  $r_0$



$$\begin{array}{r} + \boxed{\quad | \quad \boxed{\quad}} \\ - u \end{array}$$

$$\begin{array}{c} v \\ = \\ \boxed{\quad | \quad \boxed{\quad}} \\ + \boxed{\quad | \quad \boxed{\quad}} \\ + \boxed{\quad | \quad \boxed{\quad}} \\ + \boxed{\quad | \quad \boxed{\quad}} \end{array}$$

$$\begin{array}{c} u \\ = w \\ \downarrow \\ u' \\ + \\ v \\ = w \end{array}$$

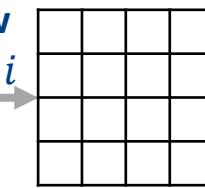
C Holds  
 $[r]_0 = r_0, [T]_0$

$(N-1, N)$ -ROT



$$M = (m_0, \dots, m_N)$$

$i$ -th row Shift by  $i$



$$\begin{array}{c} u \\ + \hat{T} \\ v \\ = [T]_1 \end{array}$$

Lookup table

Left shift by  $r_0$

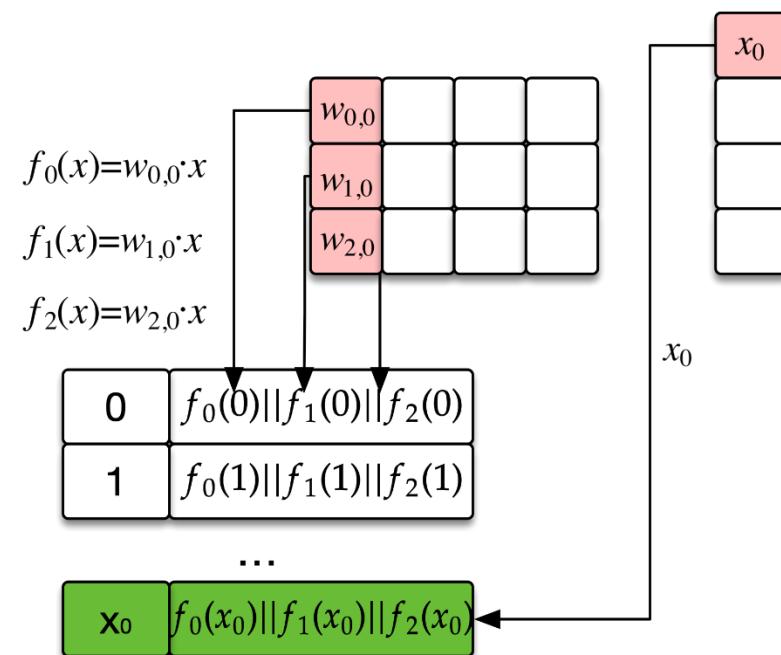
$$s' - w = [T]_0$$

S Pick  $r_1$ , Left shift  $T$  by  $r_1 = 1$  positions

S Holds  
 $[r]_1 = r_1, [T]_1$

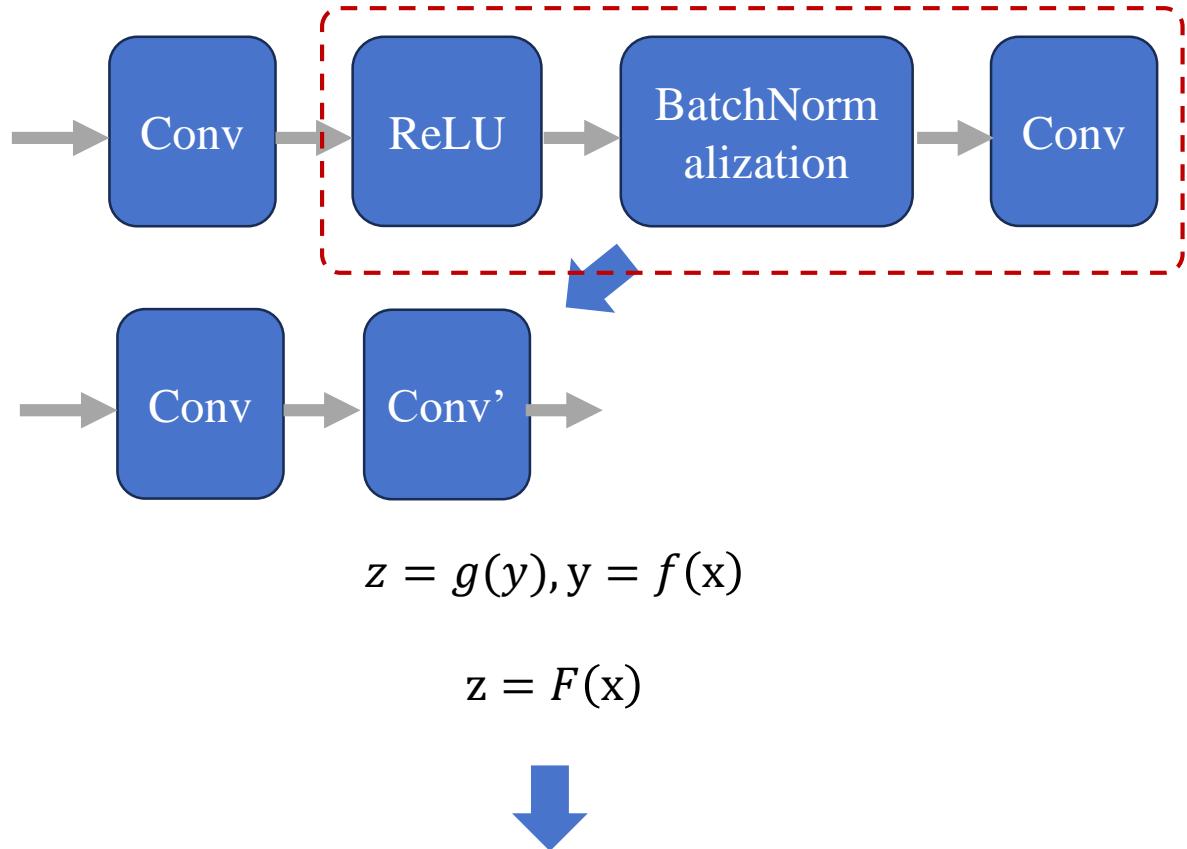
# Optimization

## Batch Scalar multiplications



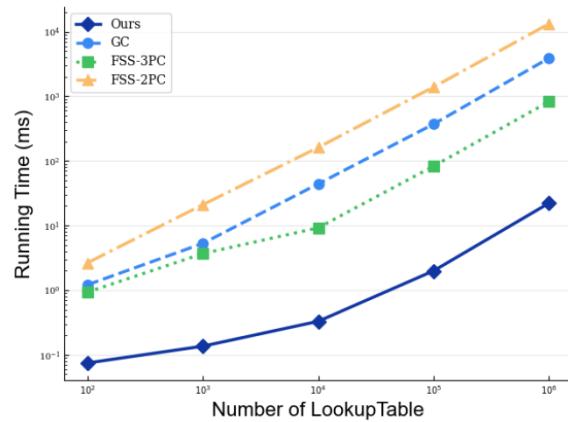
Only need to reveal  $x_0 - r$ , for all  $w_{0,i}$

## Operator Fusion

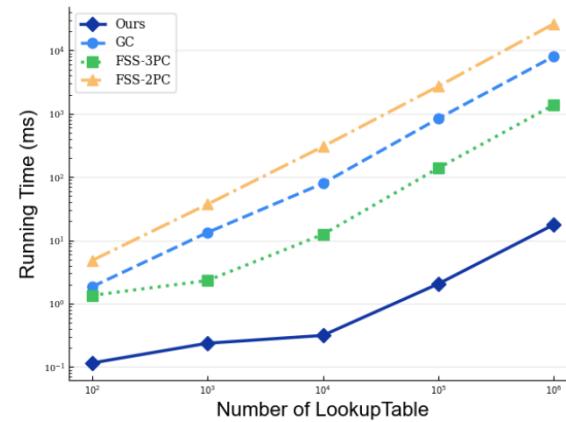


$$T := \{F(0), F(1), \dots, F(n-1)\}$$

# Performance

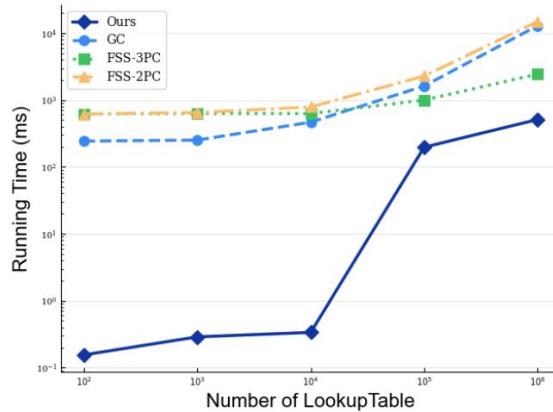


(a) 4-bit Quantization.

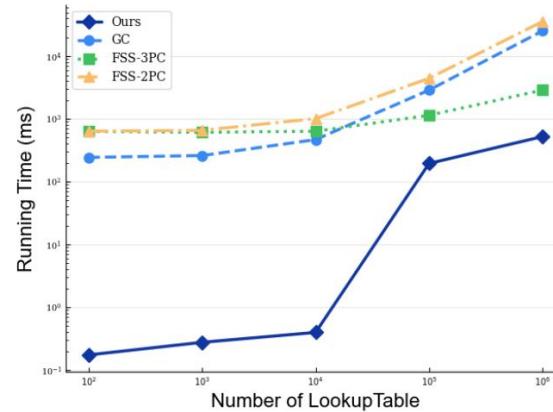


(b) 8-bit Quantization.

Online Performance comparison of Look-up table evaluation in the LAN setting



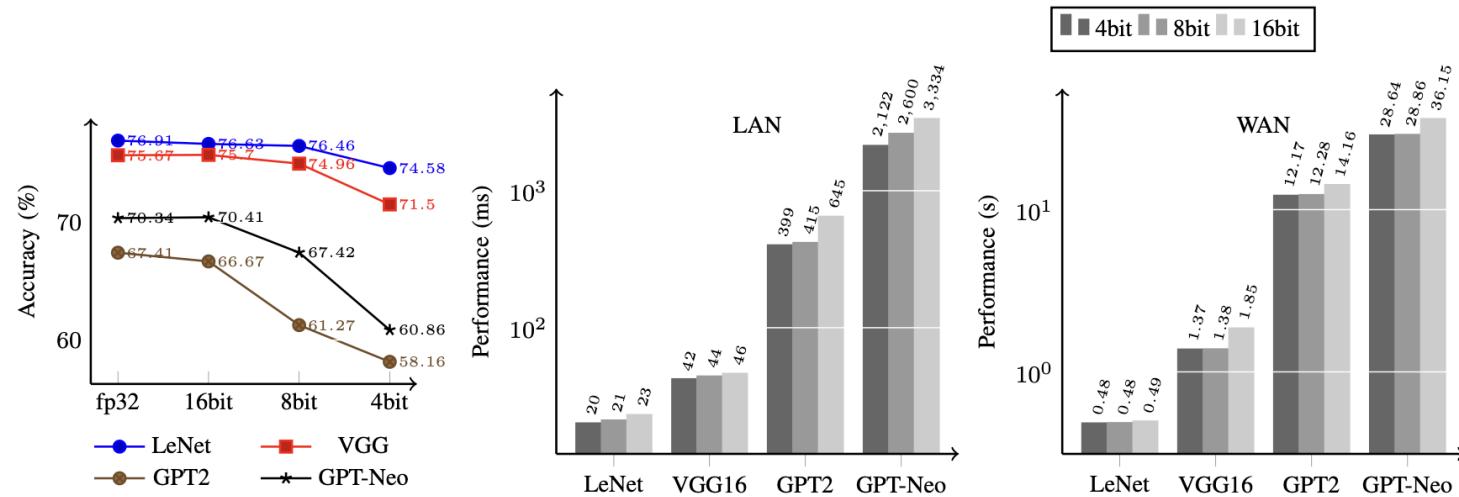
(a) 4bit Quantization.



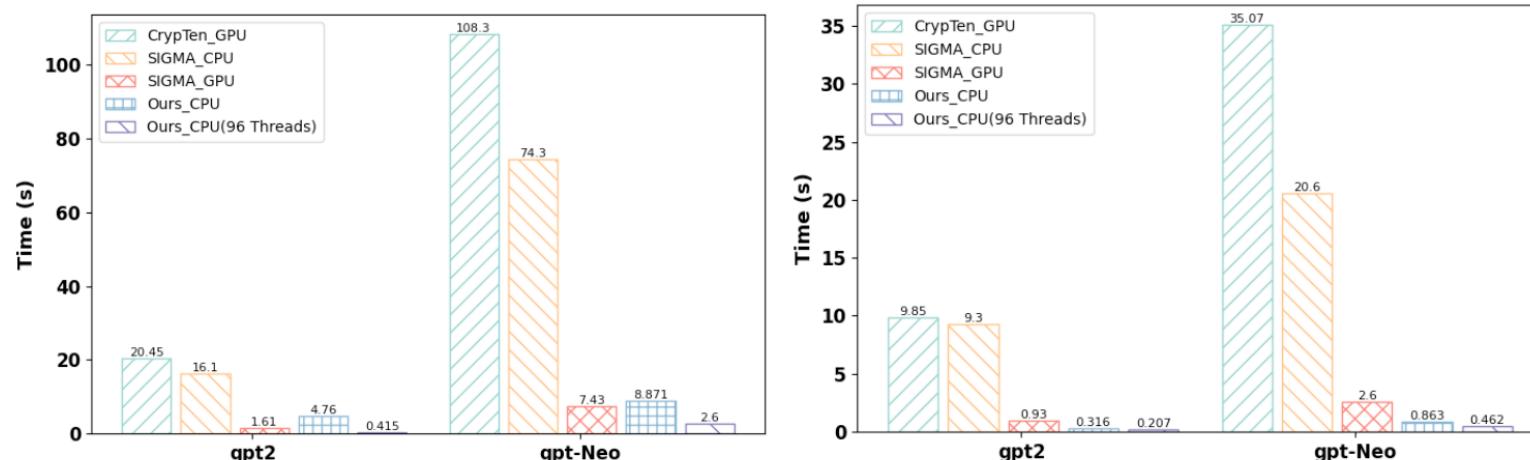
(b) 8bit Quantization.

Online Performance comparison of Look-up table evaluation in the WAN setting

# Performance



Online Performance comparison of 4-bit/8-bit/16-bit quantization



(a) LAN setting  
(b) WAN setting  
Online performance comparison for LLM with Sigma and CryptTen

- Fixed-point truncation in PPML introduces heavy overhead, hindering its use in quantized evaluations.
- Our Private Lookup Table Evaluation Scheme bypasses truncation, supporting both linear and nonlinear computations.
- Our approach achieves 20–80× speedups over non-quantized methods.



---

**Thank You**