

# Breaking the Layer Barrier: Remodeling Private Transformer Inference with Hybrid CKKS and MPC

Tianshi Xu<sup>1</sup>, Wen-jie Lu<sup>2</sup>, Jiangrui Yu<sup>1</sup>, Yi Chen<sup>1</sup>,  
Chenqi Lin<sup>1</sup>, Runsheng Wang<sup>1</sup>, Meng Li<sup>1</sup>.

*<sup>1</sup>Peking University, <sup>2</sup>TikTok.*

- Background of Private Transformer Inference
- Motivation of BLB
- Method
- Experimental Results

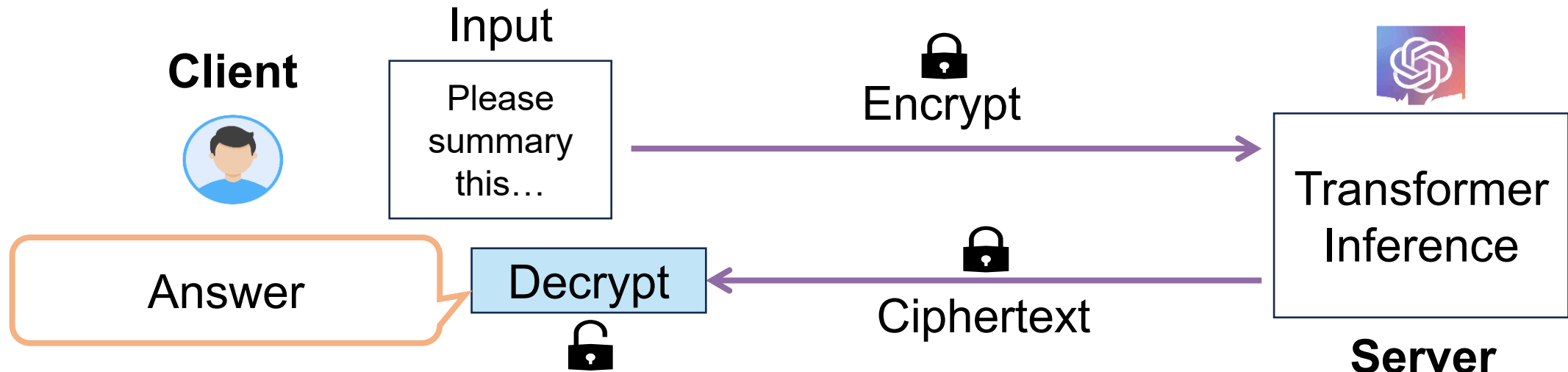
- Background of Private Transformer Inference
- Motivation of BLB
- Method
- Experimental Results

## □ Privacy-Preserving Inference for Transformer

- Transformer-based LLMs have been deployed in many **data sensitive** scenarios

## □ Privacy-Preserving Inference for Transformer

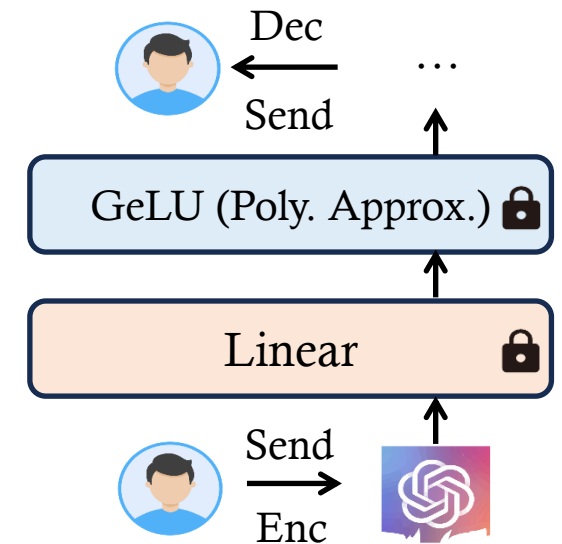
- Transformer-based LLMs have been deployed in many **data sensitive** scenarios
- **Cryptography-based** private inference offers **strong** provable protection of both **model** and **user input**



## ❑ Problems of Private Transformer Inference

### ❑ Fully Homomorphic Encryption (FHE)

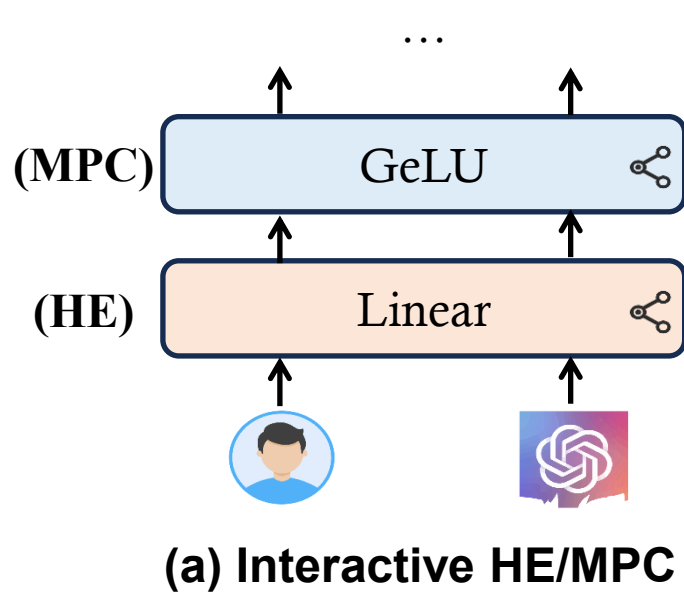
- ❑ Fail to handle nonlinear layer (**Bad accuracy** and **excessive computation** due to bootstrapping)



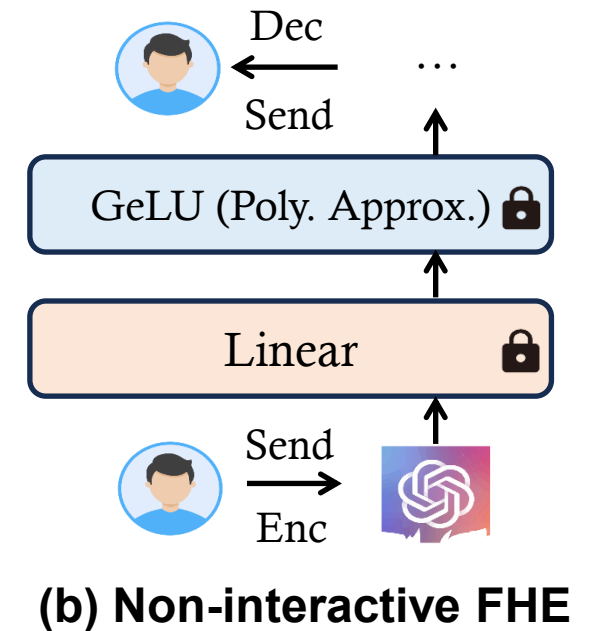
(b) Non-interactive FHE

## ❑ Problems of Private Transformer Inference

- ❑ Fully Homomorphic Encryption (FHE)
  - ❑ Fail to handle nonlinear layer (**Bad accuracy** and **excessive computation** due to bootstrapping)
- ❑ Homomorphic Encryption + Secure Multi-party Computation (hybrid **HE/MPC**)
  - ❑ High accuracy, but **excessive communication** (60GB for a BERT-base inference)

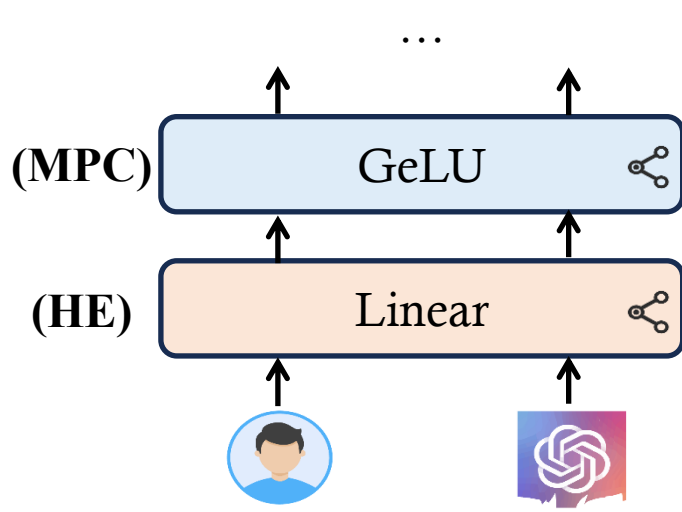


	Accuracy	Communication	Computation
HE/MPC	😊	😞	😊
FHE	😞	😊	😞
BLB (ours)	😊	😊	😊 (GPU)



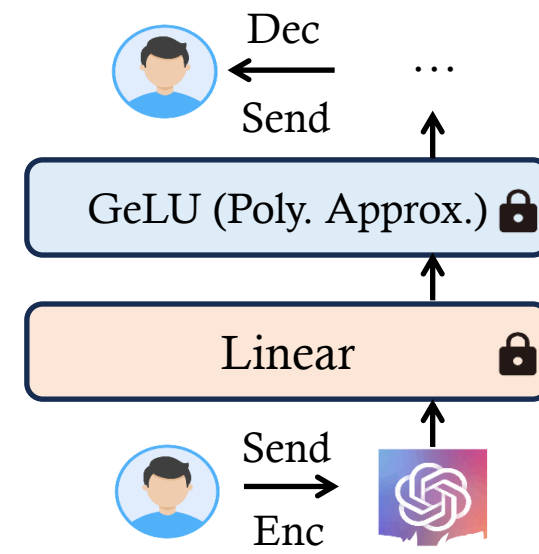
## ❑ Problems of Private Transformer Inference

- ❑ Fully Homomorphic Encryption (FHE)
  - ❑ Fail to handle nonlinear layer (**Bad accuracy** and **excessive computation** due to bootstrapping)
- ❑ Homomorphic Encryption + Secure Multi-party Computation (hybrid **HE/MPC**)
  - ❑ High accuracy, **but excessive communication** (60GB for a BERT-base inference)
- ❑ Optimization Golden Rule: **Accuracy > Communication > Computation**



(a) Interactive HE/MPC

	Accuracy	Communication	Computation
HE/MPC	😊	😞	😊
FHE	😞	😊	😞
BLB (ours)	😊	😊	😊 (GPU)

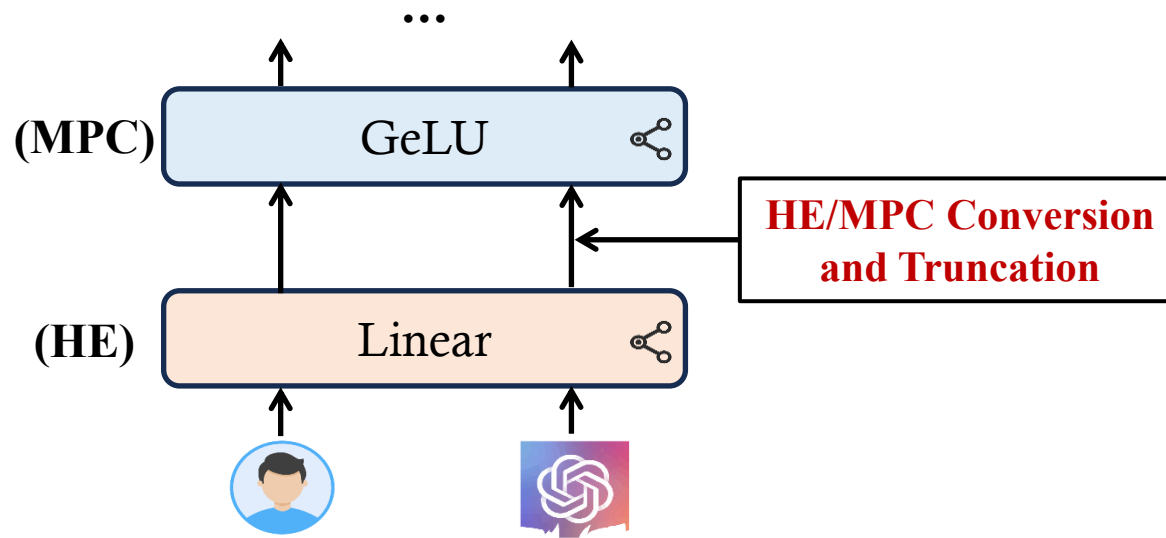


(b) Non-interactive FHE

- Background of Private Transformer Inference
- Motivation of BLB
- Method
- Experimental Results

## Private Inference with Hybrid HE/MPC

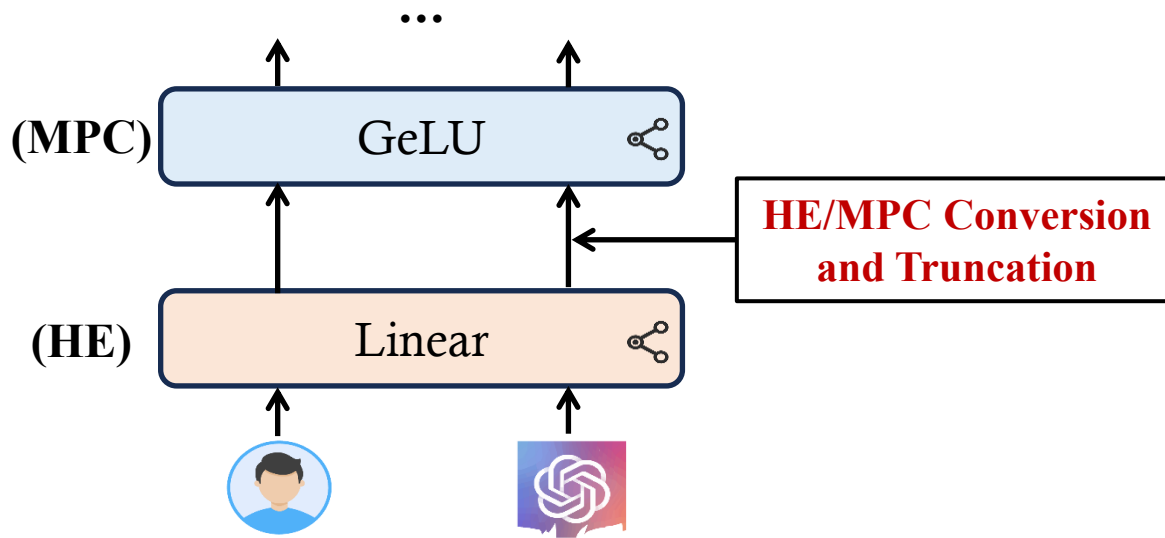
- ❑ Linear layers: Homomorphic Encryption (HE)
- ❑ Nonlinear layers: Secure Multi-party Computation (MPC)



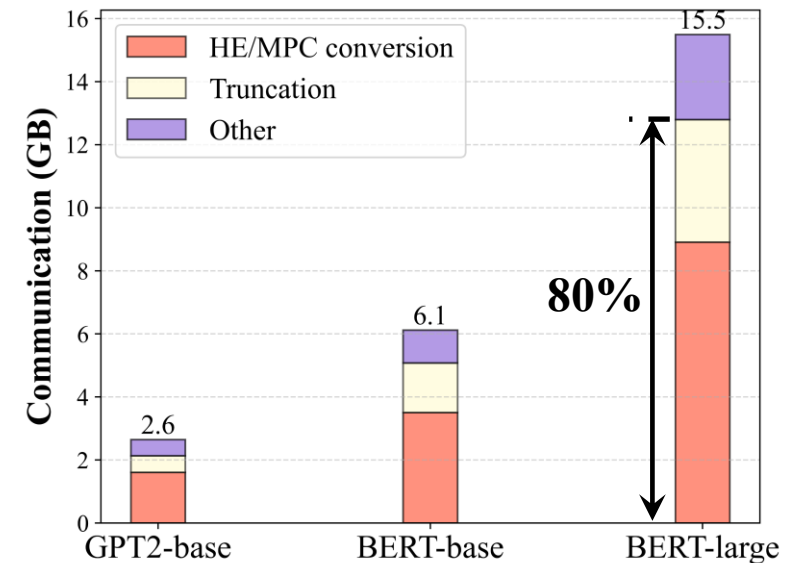
(a) Hybrid HE/MPC framework

## Private Inference with Hybrid HE/MPC

- ❑ Linear layers: Homomorphic Encryption (HE)
- ❑ Nonlinear layers: Secure Multi-party Computation (MPC)
- ❑ Communication comes from: **HE/MPC Conversions and Truncations cost dominates (>80%)**



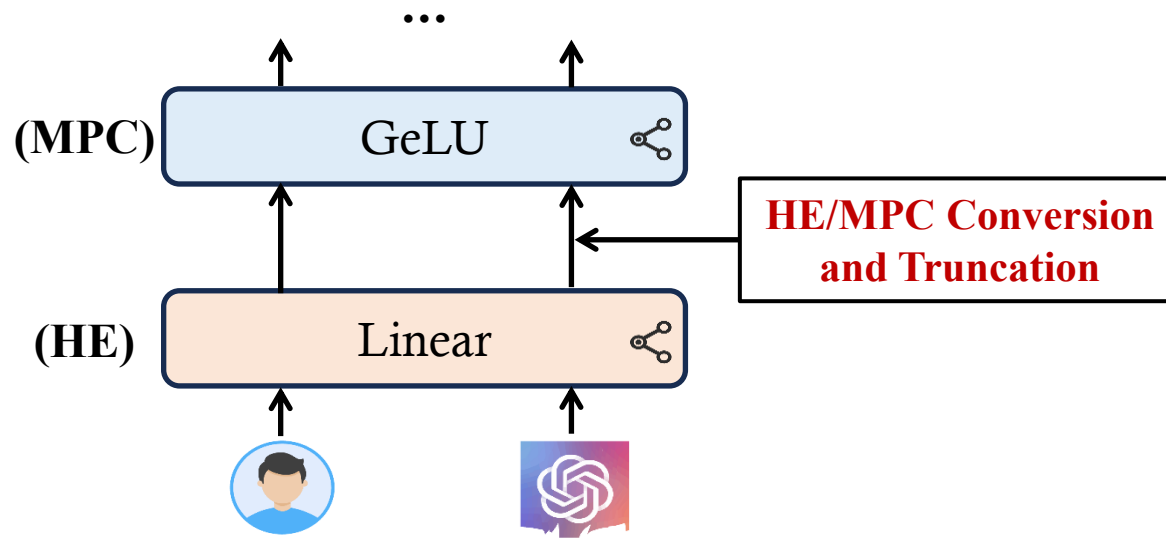
(a) Hybrid HE/MPC framework



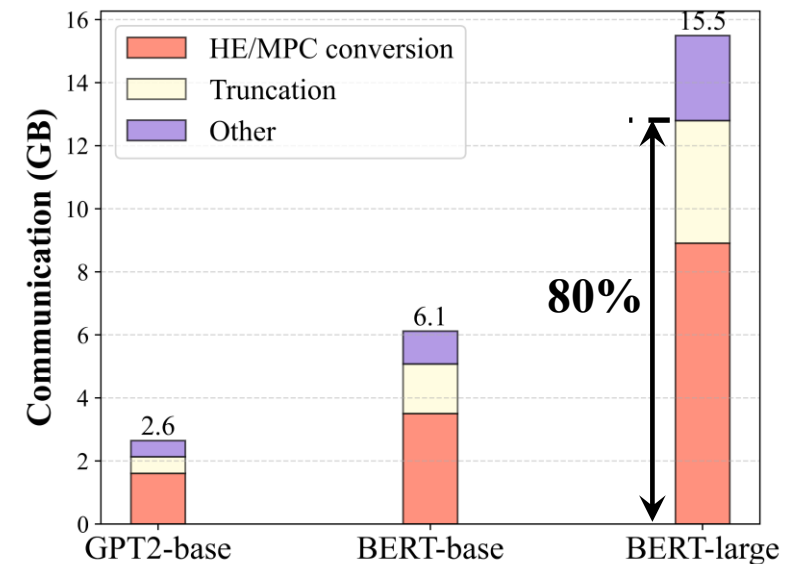
(b) Communication breakdown

## Private Inference with Hybrid HE/MPC

- ❑ Linear layers: Homomorphic Encryption (HE)
- ❑ Nonlinear layers: Secure Multi-party Computation (MPC)
- ❑ Communication comes from: **HE/MPC Conversions and Truncations cost dominates (>80%)**
- ❑ **Our Goal:** Reduce communication overhead in hybrid HE/MPC for Transformer inference



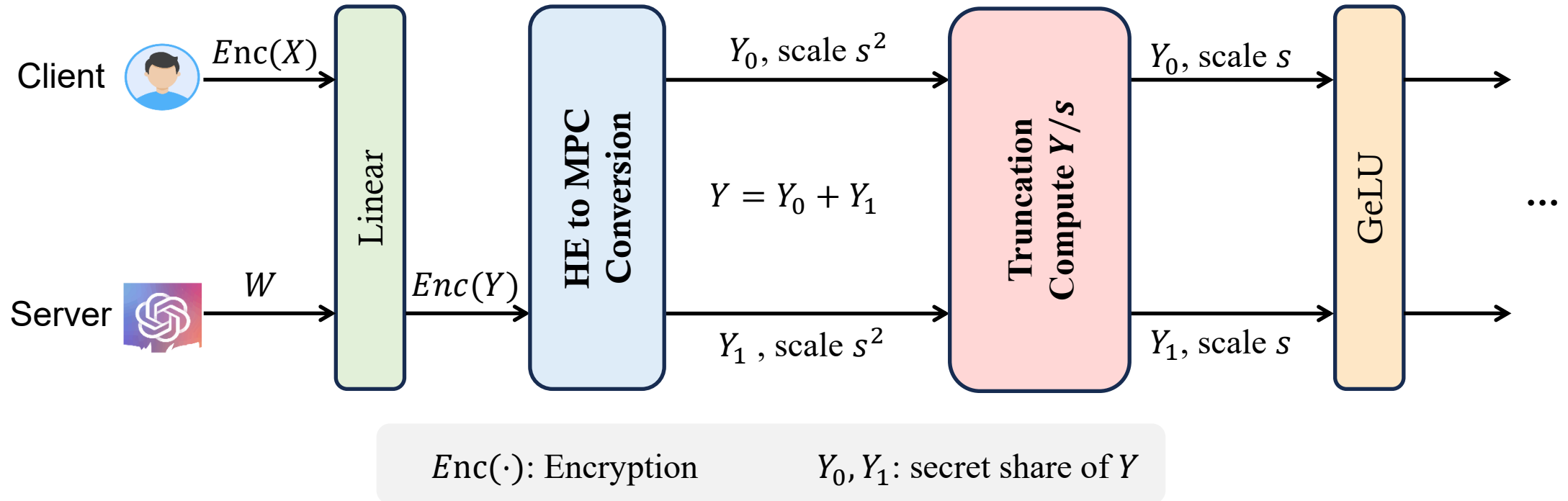
(a) Hybrid HE/MPC framework



(b) Communication breakdown

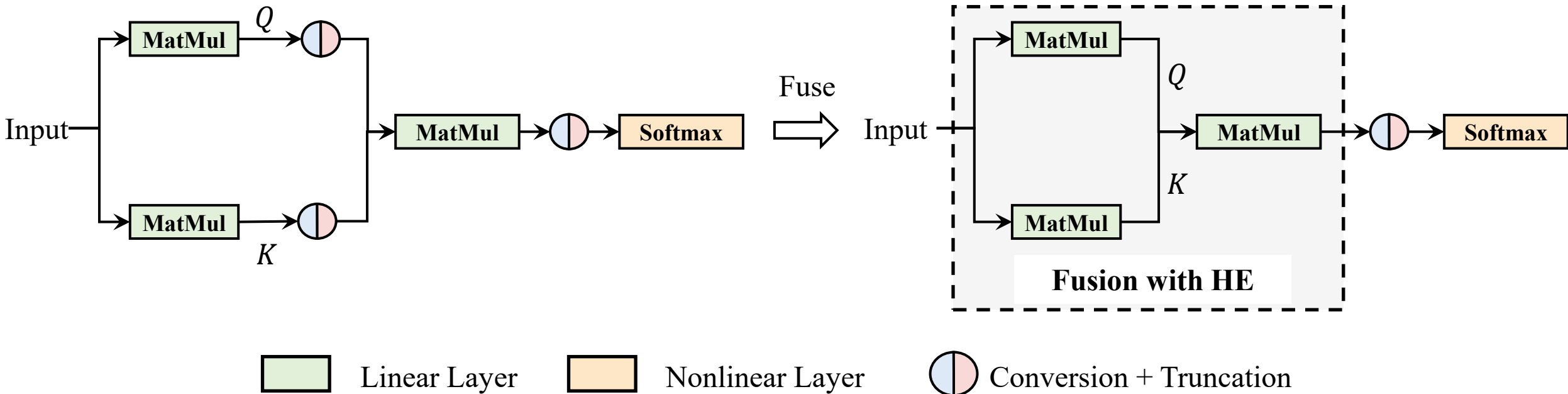
## HE/MPC Conversions and Truncations cost dominates (>80%)

- ❑ HE/MPC (Secret Share) Conversion
- ❑ Truncation: Restore scale in fixed-point arithmetic



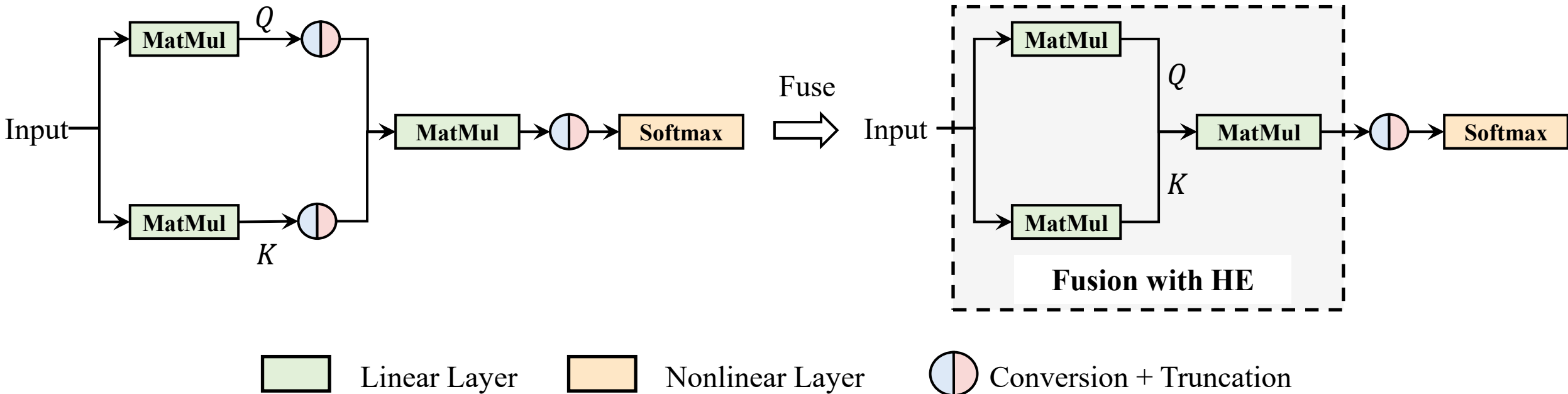
## Previous Solution: Linear Layer Fusion

- ❑ Using HE to evaluate **consecutive** linear layers **within a single communication round**
- ❑ **Eliminating all conversions and truncations** between consecutive linear layers



## Previous Solution: Linear Layer Fusion

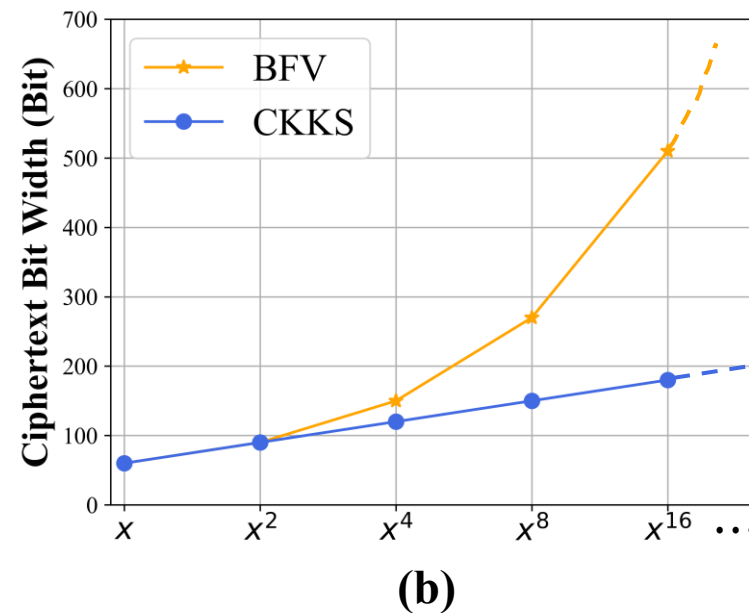
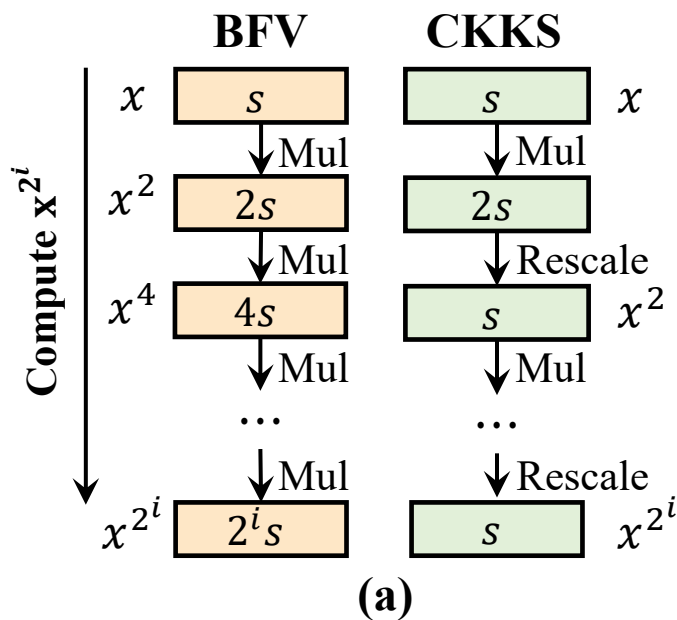
- ❑ Using HE to evaluate **consecutive** linear layers **within a single communication round**
- ❑ **Eliminating all conversions and truncations** between consecutive linear layers
- ❑ Example:  $QK^T$ ,  $Q = XW_Q$ ,  $K = XW_K$





## Three Limitations of Linear Layer Fusion

- ❑ Limitation 2: High Ciphertext Bit Width using BFV/MPC
  - ❑ BFV cannot reduce scale, needs **exponential scale growth** with multiplication depth
  - ❑ Exponential growth of scale, plaintext bit-width and ciphertext bit-width



## Three Limitations of Linear Layer Fusion

- ❑ **Limitation 3: Suboptimal MatMul Protocol after Fusion**
  - ❑ Ciphertext-Ciphertext (ct-ct) Matrix Multiplication in HE is extremely complicated!
  - ❑ Adjusting the packing of intermediate ciphertexts involves complex HE operations

## Three Limitations of Linear Layer Fusion

- ❑ **Limitation 3: Suboptimal MatMul Protocol after Fusion**
  - ❑ Ciphertext-Ciphertext (ct-ct) Matrix Multiplication in HE is extremely complicated!
  - ❑ Adjusting the packing of intermediate ciphertexts involves complex HE operations
  - ❑ BOLT(S&P'24) requires approximately **20×** more computationally-intensive **HE rotations** for fused ct-ct MatMul compared to unfused Ciphertext-Plaintext (ct-pt) protocol

- Background of Private Transformer Inference
- Motivation of BLB
- Method
- Experimental Results

**Observation:** Nonlinear layers in Transformer consist of many small linear operators

(e.g., Mul, Add)

$$\text{LayerNorm}(\mathbf{X})_{i,j} = \frac{\gamma_j(\mathbf{X}_{i,j} - \mu_i)}{\sigma_i} + \beta_j$$

$$\text{ApproxGELU}(x) = \begin{cases} x & \text{if } x > 2.7, \\ a|x|^4 + b|x|^3 + c|x|^2 & \text{if } |x| \leq 2.7, \\ + d|x| + e + 0.5x & \\ 0 & \text{if } x < -2.7. \end{cases}$$

**Observation:** Nonlinear layers in Transformer consist of many small linear operators

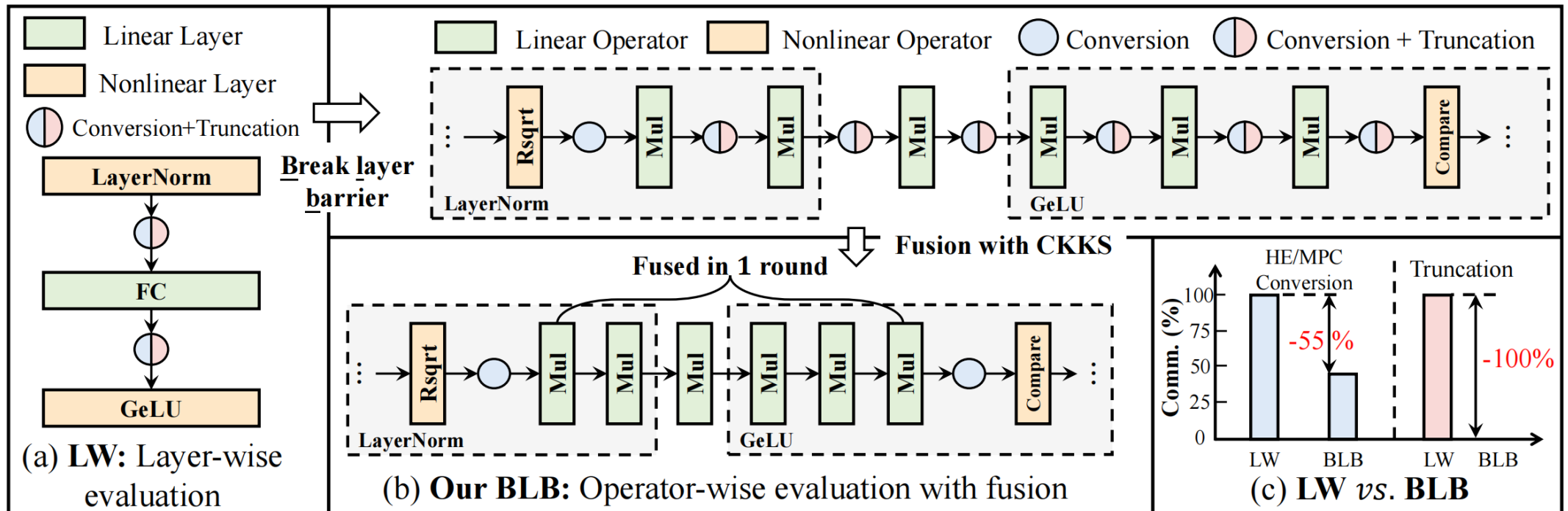
(e.g., Mul, Add)

$$\text{LayerNorm}(\mathbf{X})_{i,j} = \frac{\gamma_j(\mathbf{X}_{i,j} - \mu_i)}{\sigma_i} + \beta_j$$
$$\text{ApproxGELU}(x) = \begin{cases} x & \text{if } x > 2.7, \\ a|x|^4 + b|x|^3 + c|x|^2 & \text{if } |x| \leq 2.7, \\ + d|x| + e + 0.5x & \\ 0 & \text{if } x < -2.7. \end{cases}$$

**Intuition:** Nonlinear layers should **no** longer be treated as indivisible units. A finer-grained analysis is required.

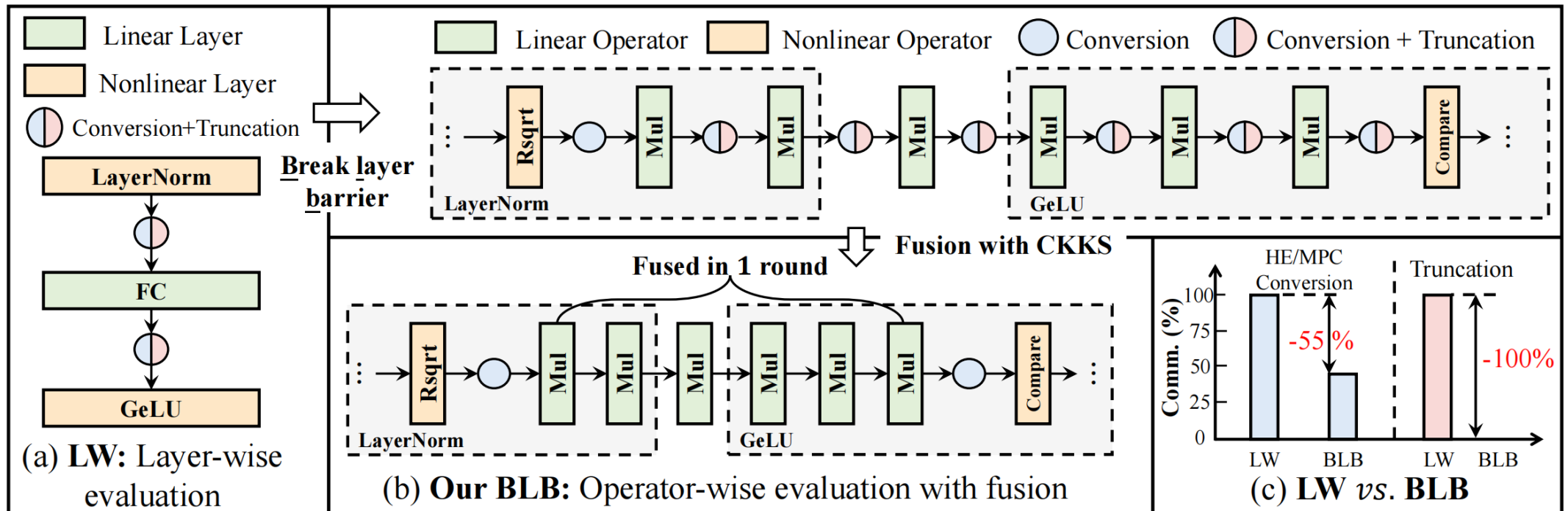
## BLB: From **layer-wise (LW)** evaluation to **operator-wise evaluation with fusion**

- Decompose layers into operators, classifying them into linear operators and nonlinear operators



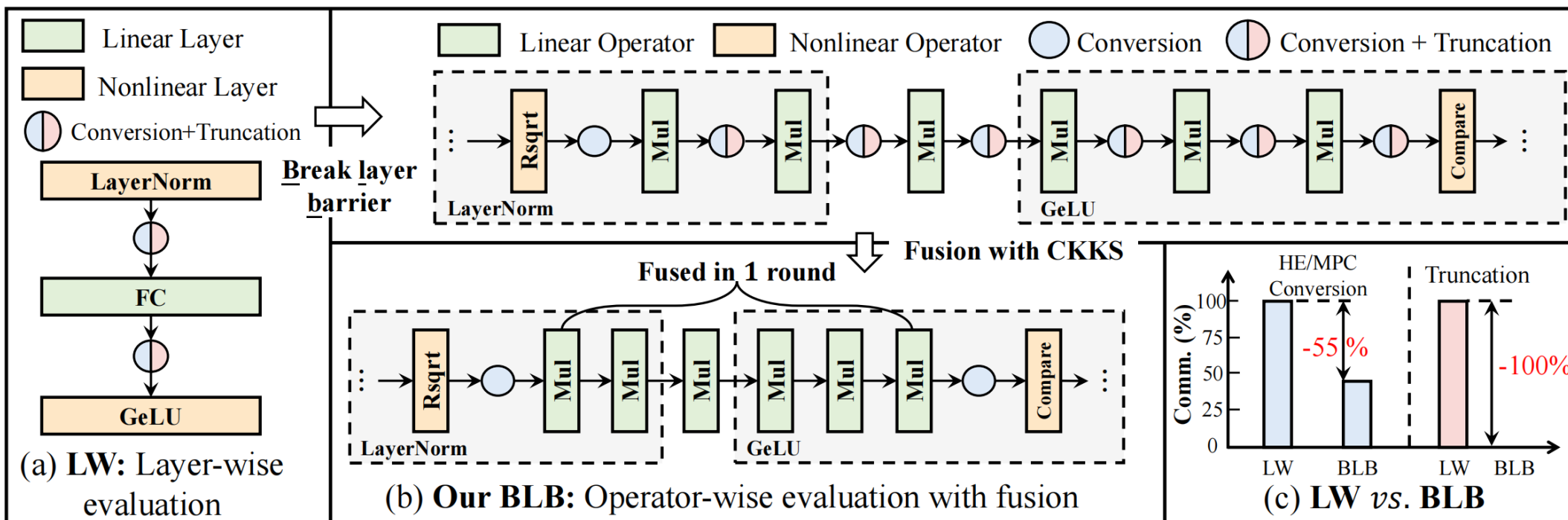
## BLB: From **layer-wise (LW)** evaluation to **operator-wise evaluation with fusion**

- ❑ Decompose layers into operators, classifying them into linear operators and nonlinear operators
- ❑ Fuse adjacent linear operators across layers (evaluate them within a communication round)



## BLB: From **layer-wise (LW)** evaluation to **operator-wise evaluation with fusion**

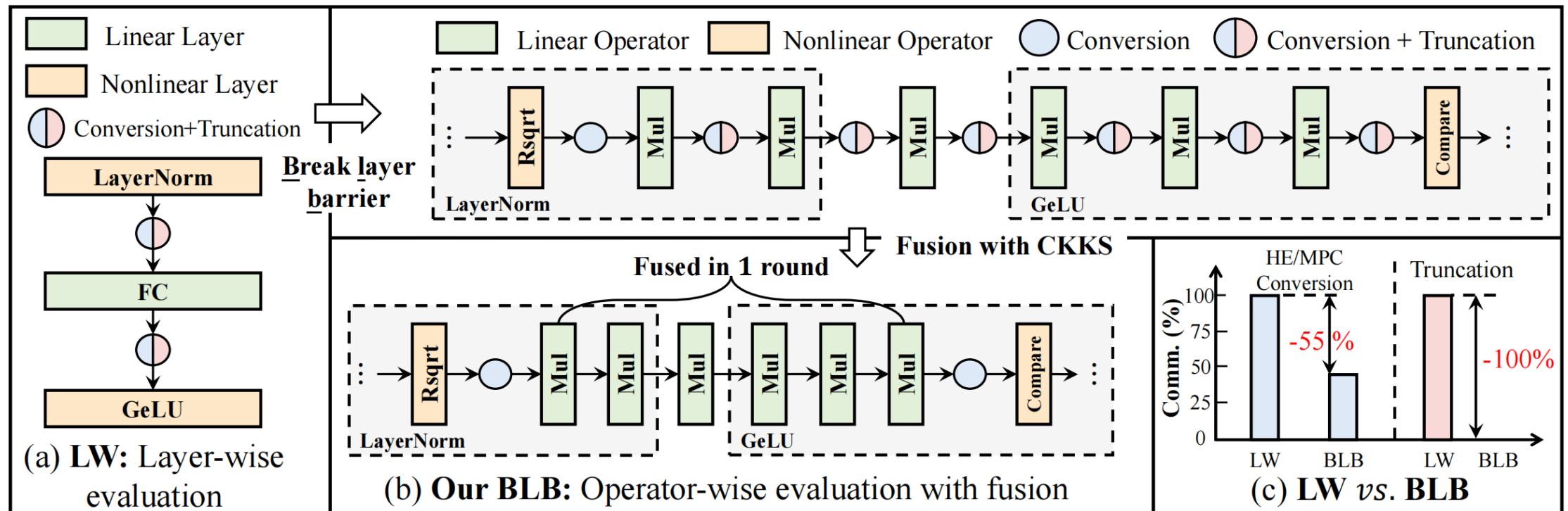
- ❑ Decompose layers into operators, classifying them into linear operators and nonlinear operators
- ❑ Fuse adjacent linear operators across layers (evaluate them within a communication round)
- ❑ Eliminate all conversions and truncations between adjacent linear operators



# Core Idea: Breaking the Layer Barrier

**Example:** LayerNorm + Fully Connected Layer (FC) + GeLU

- ❑ The last two ops in LayerNorm, the subsequent FC Layer, the first three ops in GeLU can be fused
- ❑ **100% Truncation** and **55% HE/MPC Conversion** can be eliminated after fusion



## Three Technical Highlights

### 1. FineGrainFusion: Systematic operator-level fusion patterns across Transformer layers

Type	Name	Description
<i>Linear operators</i>		
Identity	ewadd <sub>cc</sub> , ewadd <sub>cp</sub>	Element-wise addition of ct-ct, ct-pt
	ewmul <sub>cc</sub> , ewmul <sub>cp</sub>	Element-wise multiplication of ct-ct, ct-pt
Expansion	sadd <sub>cc</sub> , sadd <sub>cp</sub>	Scalar addition of ct-ct, ct-pt
	smul <sub>cc</sub> , smul <sub>cp</sub>	Scalar multiplication of ct-ct, ct-pt
Reduction	sum	Sum in a specified dimension of the input
Transformation	matmul <sub>cc</sub> , matmul <sub>cp</sub>	MatMul of ct-ct, ct-pt
<i>Nonlinear operators</i>		
	cmp	$\llbracket \mathbf{1}\{x < y\} \rrbracket^B \leftarrow \text{cmp}(\llbracket x \rrbracket^M, \llbracket y \rrbracket^M)$
	mux	$\llbracket b \cdot x \rrbracket^M \leftarrow \text{mux}(\llbracket b \rrbracket^B, \llbracket x \rrbracket^M)$
	rec, rsqrt	Reciprocal, Reciprocal Sqrt [47]

### Operator Categorization

Second op / First op	Identity	Expansion	Reduction	Transformation
Identity	Identity	Expansion	Reduction	Transformation
Expansion	Expansion	×	Identity	Expansion
Reduction	Reduction	Identity	×	×
Transformation	Transformation	×	Reduction	Transformation

### Fusion pattern analysis for linear operators

## Three Technical Highlights

### 2. CKKS+MPC Conversion: **First secure** protocol to support CKKS/MPC conversion

- We use **CKKS** instead of BFV to control the ciphertext bitwidth growth
- Previous CKKS-MPC conversion protocol is insecure, leaking computation results

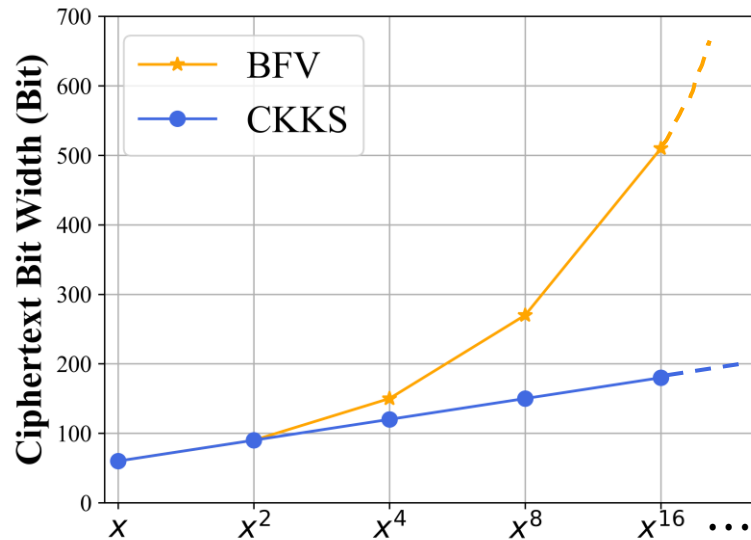


Fig: CKKS ct biwdith grows linearly

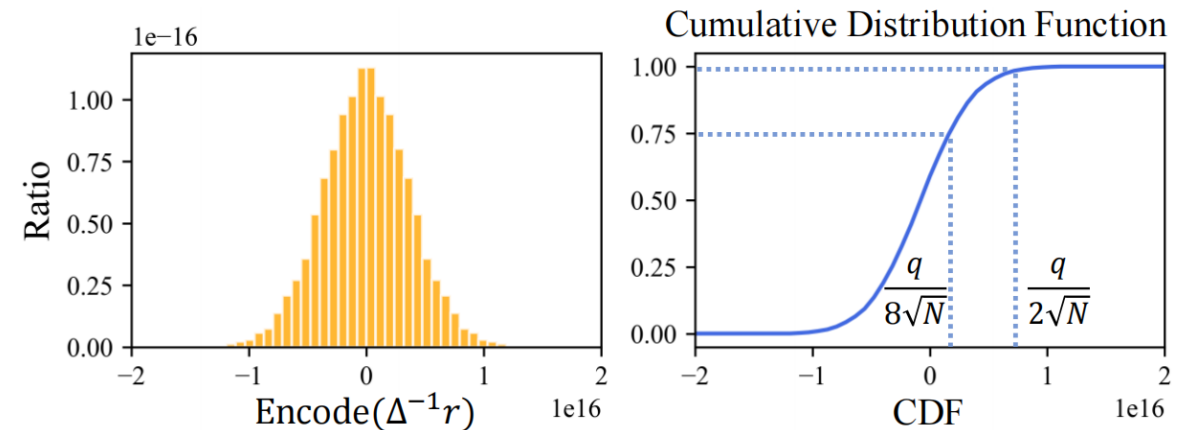


Fig: CKKS Encode function narrows random mask

## Three Technical Highlights

3. **Rotation-Efficient Ct-Ct MatMul:** 29 × and 8 × fewer HE rotations than BOLT (S&P'24) and Powerformer (ACL'25)
- Optimized for  $QK^T$  and Softmax ×  $V$ ; multi-head packing; BSGS optimization

matmul <sub>cc</sub>		BOLT [45]	Powerformer [46]	BLB
$Q_h K_h^T$ $h \in [H]$	# Rot	$O(\frac{md}{s}(m + \log_2 m))$ 18432	$O(\frac{md}{s}(\frac{5}{2}m + 4\sqrt{m}))$ 5856	$O(\frac{md}{s}(\frac{\sqrt{d}}{H} + \sqrt{m}) + m)$ 640
	# CMult.	$O(\frac{m^2 d}{s})$ 2048	$O(\frac{m^2 d}{s})$ 1024	$O(\frac{m^2 d}{s})$ 1024
Softmax × $V_h$ $h \in [H]$	# Rot	$O(\frac{m^2 H}{s}(m + \log_2 m))$ 28560	$O(\frac{m^2 H}{s}(\frac{5}{2}m + 4\sqrt{m}))$ 6508	$O(\frac{m^2 H}{s}(\sqrt{\frac{m}{H}} + \sqrt{m}) + m)$ 1056
	# CMult.	$O(\frac{m^2 d}{s})$ 1024	$O(\frac{m^3 H}{s})$ 2048	$O(\frac{m^3 H}{s})$ 2048

- Background of Private Transformer Inference
- Motivation of BLB
- Method
- Experimental Results

## Experimental Setup

- ❑ Implementation: SEAL, EzPC, Phantom FHE (GPU)
- ❑ Models: GPT2-base, BERT-base, BERT-large
- ❑ Baselines
  - ❑ HE/2PC: Iron (NeurIPS'22), BOLT (S&P'24), Bumblebee (NDSS'25)
  - ❑ 3PC: SIGMA (PETS'24), MPCFormer (ICLR'23)
  - ❑ FHE: NEXUS (NDSS'25)

## Communication Reduction

- **21 ×** vs. BOLT (S&P'24)
- **2 ×** vs. Bumblebee (NDSS'25)

## Latency Reduction

- Up to **29 ×** on CPU
- Up to **13 ×** on GPU

Model	Framework	Latency (min)			Comm. (GB)
		LAN	WAN <sub>2</sub>	WAN <sub>3</sub>	
GPT2-base 64 input tokens	SIGMA [26]*	7.7	25.1	38.5	28.7
	MPCFormer [39]*	2.1	7.2	10.5	7.3
	BOLT [45]	9.7	34.0	53.6	34.8
	Bumblebee [41]	<b>2.9</b>	<b>5.2</b>	12.3	2.5
	BLB	4.0	6.8	<b>10.2</b>	<b>1.5</b>
	BOLT (GPU) [45]	7.6	31.8	51.6	34.8
	Bumblebee (GPU) [41]	2.6	4.9	12.0	2.5
	BLB (GPU)	<b>2.0</b>	<b>3.9</b>	<b>8.1</b>	<b>1.5</b>
	BERT-base 128 input tokens	NEXUS [56]	457	458	470
	Iron [27]	66.7	/	/	76.5
	SIGMA [26]*	7.8	30.4	48.6	34.4
	MPCFormer [39]*	5.5	18.0	27.5	12.1
	BOLT [45]	22.1	85.0	130.0	63.6
	Bumblebee [41]	<b>5.0</b>	12.3	22.1	5.8
	BLB	6.1	<b>10.5</b>	<b>17.2</b>	<b>3.0</b>
	NEXUS (GPU) [56]	19.9	20.8	32.5	0.2
	BOLT (GPU) [45]	15.5	77.9	122.7	63.6
	Bumblebee (GPU) [41]	4.3	11.6	21.3	5.8
	BLB (GPU)	<b>2.5</b>	<b>6.6</b>	<b>13.2</b>	<b>3.0</b>
BERT-large 128 input tokens	Iron [27]	92.0	/	/	220.0
	SIGMA [26]*	20.5	102.5	155.5	92.8
	MPCFormer [39]*	7.7	34.1	52.0	32.6
	BOLT [45]	57.6	222.0	274.4	158.9
	Bumblebee [41]	<b>10.6</b>	32.4	51.3	15.2
	BLB	15.1	<b>29.0</b>	<b>37.7</b>	<b>7.8</b>
	BOLT (GPU) [45]	43.8	208.2	260.6	158.9
	Bumblebee (GPU) [41]	9.7	30.8	49.2	15.2
	BLB (GPU)	<b>6.6</b>	<b>16.2</b>	<b>24.9</b>	<b>7.8</b>

\* Frameworks that require three parties.

- ❑ BLB breaks the “**one-layer, one-protocol**” paradigm
- ❑ Enables **operator-level fusion** with CKKS+MPC
- ❑ Achieves **significant communication and latency reduction**

*From layers to operators — it's time to break the barrier*