Secure Transformer Inference Made Non-interactive

Jiawen Zhang*, Xinpeng Yang*, Lipeng He‡, Kejia Chen*, Wen-jie Lu*, Yinghao Wang*, Xiaoyang Hou*, Jian Liu*, Kui Ren*, Xiaohu Yang*

*Zhejiang University, *‡University* of Waterloo





Secure ML Inference



Transformer vs CNN Inference

In the context of non-interactive secure inference, there are two key differences between transformers and CNN models:

Larger Scale Matrix-Matrix Multiplications

Gazalle [USENIX Security '18], Cheetah [USENIX Security '22] and Iron compute matrix multiplications via **inner dot products**, and employ **sparse packing** for the resulting ciphertexts. The majority of the data slots in the ciphertexts are wasted.

Input dimensions of matrices in transformer-based models are much higher, leading to **increased computational cost**.

Higher Dimensional Inputs to Argmax

Input to Argmax is a probability vector of size m. In CNN classification tasks, m is usually not large (m = 1,000 in ImageNet-1k). However, in transformer-based NLP tasks, m = vocabulary size, with m reaching 30,522 in BERT and 128,256 in Llama-3-8B.

Previous SOTA Phoenix [CCS '22] can only achieves a computational complexity of O(m)

Prior Work: Interactive Protocols

Interactive Secure Transformer Inference Protocols

- Iron [NeurIPS '22]
- BOLT [IEEE S&P '24]
- BumbleBee [NDSS '25]

Based on Multi-party Secure Computation (MPC)

- $\circ~$ Substantial communication overhead
- Unable to support hardware acceleration via GPU, FPGA, etc.

e.g. BOLT requires:

- \circ 59.61 GB of bandwidth
- **10,509** interaction rounds
- **\$5.44** per response token

Our Work: NEXUS

Based on RNS-CKKS Fully Homomorphic Encryption (FHE), we achieve:



Efficient and Communicationoptimized **Matrix Multiplication**



Efficient **Argmax** and other **nonlinear function** evaluation



CPU/GPU Implementation

Efficient Matrix Multiplication Toy Example

 $\pmb{A} \in \mathbb{R}^{2 imes 3}, \pmb{W}_Q \in \mathbb{R}^{3 imes 3},$



Efficient Matrix Multiplication Secure ciphertext compression



Efficient Matrix Multiplication Secure ciphertext decompression [SealPIR S&P'18]



Efficient Matrix Multiplication Results Comparison

BERT-base: m = 256, n = 768, k = 64Amortized cost of t = 256 matrix multiplications

Methods	# Cipher	texts	# Key switching		
Iron [NeurIPS '22]	$rac{2\sqrt{mnk}}{\sqrt{N}}$	111	-	0	
Bumblebee [NDSS '25]	$rac{m(n+k)}{N}$	52	$\frac{mk\log(N)}{2\sqrt{N}}$	1536	
BOLT [IEEE S&P '24]	$rac{m(n+k)}{N}$	52	$\sqrt{rac{m^2n^2k}{N^2}}$	384	
Ours	$\frac{mk}{N} + \frac{nk}{Nt}$	5	$rac{2nk-2}{t}$	384	



⁽b) Amortized Communication vs. #inputs.

Efficient Argmax evaluation Preliminaries

Argmax: Returns the indices of the maximum values along an vector. The Argmax of [2, 3, 4, 1] is [0, 0, 1, 0]



Efficient **Argmax** evaluation



$y = Sign(x - x_{max}) + 1$								
x	0	0	0	0	1	4	3	2
y = x + RotR(x, 4)	1	4	3	2	1	4	3	2
$r = RotL(y, 2^0)$	#	1	4	3	2	1	4	3

- $tL(y, 2^0)$
- y = Max(r, y)3
- $r = RotL(y, 2^1)$
- y = Max(r, y)
 - y = Mask(y)

Ours: O(log N)

Efficient **Argmax** evaluation



Computing the root of a *binary tree*: $O(\log N)$ *f*() can be any associative, such as sum, max, min...

Other Non-linear Function evaluation GELU, Softmax, LayerNorm

Following similar methods presented in BumbleBee [NDSS '25] and PUMA [DLZ+23]:

GELU: approximate using a piecewise polynomial:

$$\operatorname{GELU}(x) = \begin{cases} 0 & x \leq -4 \\ P(x) = \sum_{i=0}^{i=3} c_i x^i & -4 < x \leq -1.95 \\ Q(x) = \sum_{i=0}^{i=6} d_i x^i & -1.95 < x \leq 3 \\ x & x > 3 \end{cases}$$

Softmax: taking a_{max} as a constant as its value does not affect the result of Softmax:

 $y_{i} = \frac{\exp(a_{i} - a_{max})}{\sum_{j=0}^{m-1} \exp(a_{j} - a_{max})} \longrightarrow \text{Goldschmidt division algorithm}$ QuickSum $EXP(x) \approx (1 + \frac{x}{2^{r}})^{2^{r}}, \quad x \leq 0$ Approximated using this Taylor series

LayerNorm:

$$y_{i} = \gamma \cdot \frac{a_{i} - \mu}{\sigma} + \beta$$

$$= \gamma \cdot \frac{n(a_{i} - \mu)}{n\sqrt{\frac{1}{n}\sum_{i=0}^{n-1}(a_{i} - \mu)^{2}}} + \beta$$

$$= \sqrt{n}\gamma \cdot \frac{na_{i} - n\mu}{\sqrt{\sum_{i=0}^{n-1}(na_{i} - n\mu)^{2}}} + \beta.$$
Let $z_{i} = na_{i} - n\mu = na_{i} - \sum_{i=0}^{n-1}a_{i}$, then
$$y_{i} = \gamma\sqrt{n} \cdot \frac{z_{i}}{\sqrt{\sum_{i=0}^{n-1}z_{i}^{2}}} + \beta.$$

- QuickSum: $f(\tilde{a}, \tilde{b})$ is given by $\tilde{a} \boxplus \tilde{b}$
- Inverse Square Root: Newton's iteration

Efficient Non-linear Function Eval Results Comparison



* GPU accelerated

Placement of Bootstrapping



End-to-end Results

Onenetion	Depth	BERT-base (12 layers)			
Operation		Input	CPU(s)	GPU(s)	
MATRIXMUL	$21 \rightarrow 20$	$(\mathbb{R}^{128\times768}\times\mathbb{R}^{768\times768})\times3$	65	2.68	
MATRIXMUL	$20 \rightarrow 21$	$(\mathbb{R}^{128\times 64} \times \mathbb{R}^{64\times 128}) \times 12$	14	0.54	
Softmax	$19 \rightarrow 3$	$(\mathbb{R}^{128\times 128})\times 12$	47	1.15	
MATRIXMUL	$3 \rightarrow 2$	$(\mathbb{R}^{128\times128}\times\mathbb{R}^{128\times64})\times12$	9	0.36	
MATRIXMUL	$2 \rightarrow 1$	$\mathbb{R}^{128\times768}\times\mathbb{R}^{768\times768}$	2	0.06	
BOOTSTRAPPING	$1 \rightarrow 17$	$\mathbb{R}^{128 imes 768}$	127	5.63	
LAYERNORM	$17 \rightarrow 1$	$\mathbb{R}^{128 imes 768}$	16	1.01	
BOOTSTRAPPING	$1 \rightarrow 17$	$\mathbb{R}^{128 imes 768}$	127	5.63	
MATRIXMUL	$17 \rightarrow 16$	$\mathbb{R}^{128\times768}\times\mathbb{R}^{768\times3072}$	48	1.71	
GELU	$16 \rightarrow 2$	$\mathbb{R}^{128 imes 3072}$	44	3.35	
MATRIXMUL	$2 \rightarrow 1$	$\mathbb{R}^{128\times3072}\times\mathbb{R}^{3072\times768}$	8	0.20	
BOOTSTRAPPING	$1 \rightarrow 17$	$\mathbb{R}^{128 imes 768}$	127	5.63	
LAYERNORM	$17 \rightarrow 1$	$\mathbb{R}^{128 imes 768}$	16	1.01	
BOOTSTRAPPING	$1 \rightarrow 21$	$\mathbb{R}^{128 imes 768}$	153	5.90	
ARGMAX	*	\mathbb{R}^{30522}	54	2.48	
Total	-	-	857	37.34	

End-to-end Results



Code

📃 🌍 zju-abclab / NEXUS		Q Type 🕖 to search	🛯 🖷 • + • 💿 🟗 🖨 🌾		
<> Code Issues 4 It Pull requests	s 🕞 Actions 🖽 Projects	③ Security 🗠 Insights			
	Ŕ	Edit Pins 👻 ⓒ Unwatch 2 👻	💱 Fork 8 → 📌 Starred 86 →		
រះ main 👻 រិ Branch 🛇 0 Tags	Q Go to file	t + <> Code -	About		
💽 ttttonyhe chore: update readme.md	e150580	c · 5 months ago 🕚 110 Commits	Non-interactive protocol for secure transformer inference based on RNS- CKKS		
🖿 cuda	chore: update readme	6 months ago			
🖿 data	feat: cpu argmax	7 months ago	中 GPL-3.0 license		
Src 🖿	chore: update readme	6 months ago	Cite this repository -		
thirdparty/SEAL-4.1-bs	chore: clean up	6 months ago	-∿ Activity		
🗋 .gitattributes	chore: clean up	7 months ago	☆ 86 stars		
🗋 .gitignore	chore: update readme	7 months ago	 		
CITATION.bib	chore: add citation.bib	5 months ago			
CMakeLists.txt	chore: clean up	6 months ago			
	feat: add a license	7 months ago	Releases		
README.md	chore: update readme.md	5 months ago	No releases published Create a new release		
C README A GPL-3.0 license		Ø ∷≣	Packages		
			No packages published Publish your first package		
NEXUS			Contributors 3		
NEXUS is the first non-interactive pro	Ttttonyhe Tony L. He				
			Kevin-Zh-CS Kevin Zhang		

Summary

- We propose NEXUS, the first <u>non-interactive</u> secure <u>transformer</u> inference protocol.
- Based on RNS-CKKS Fully Homomorphic Encryption scheme, we design a series of efficient and communication-optimized matrix multiplication and non-linear function evaluation protocols.
- NEXUS **reduces bandwidth consumption by 372**. **5** × compared to BOLT [Oakland '24] and 53.6 × compared to BumbleBee [NDSS '25]. Its non-interactive property allows for better hardware acceleration, with the GPU version achieving a **42**. **3** × **speedup** in runtime.



