

---

# Characterizing and Improving MPC-based Private Inference for Transformer-based Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Secure multi-party computation (MPC) is gaining popularity with the growing  
2 demand for privacy-preserving cloud services. While there has been plenty of  
3 attention paid to MPCs for convolution neural networks (CNNs) [1, 2, 3, 4, 5],  
4 MPC-based private inference for Transformer models has not been studied in  
5 detail. This paper provides a characterization study of the performance overhead  
6 for running Transformer models with secure MPC, and proposes an optimization  
7 for embedding tables. Our study shows that Transformers introduce a couple of  
8 new challenges for MPC-based private inference: softmax and embedded tables. To  
9 address the overhead of embedding table accesses under MPC, we propose to use  
10 tensor-train (TT) decomposition, a mechanism that splits a large embedding tables  
11 into multiple smaller embedding tables. For the NLP workloads, the experiments  
12 show that the TT decomposition can speed up embedding table accesses by **2x** with  
13 only a 1.19 drop in the masked-language model perplexity score.

## 14 1 Introduction

15 Data privacy is a pressing concern for privacy-preserving machine learning (PPML) in the cloud. To  
16 address this challenge, secure computation techniques such as homomorphic encryption (HE) [6],  
17 trusted execution environments (TEE) [7, 8, 9, 10, 11, 12], and secure multi-party computation  
18 (MPC) [13] have been applied to PPML. The previous studies, however, focused primarily on  
19 convolutional neural networks (CNNs) [1, 2, 3, 4, 5]. This paper studies MPC-based private inference  
20 for Transformer-based models, which are commonly used for natural language processing (NLP) and  
21 have been applied to computer vision [14] more recently.

22 As the first step to enable efficient MPC-based private inference for Transformer models, we per-  
23 formed a detailed study of the performance overhead. The study uncovers two challenges, which  
24 are new to Transformers and have not been studied in the context of CNNs: softmax and embedding  
25 tables. While non-linear activation functions such as ReLU are known to be a major source of  
26 performance overhead for running CNNs in MPC, softmax accounts for an even larger portion of the  
27 MPC execution time for Transformers. Unlike CNNs where softmax only needs to be performed at  
28 the end, Transformers use softmax in each layer. An in-depth investigation also shows that the max  
29 function that is used for numerical stability is the main source of softmax overhead.

30 Another new challenge for Transforms comes from embedding tables, which convert categorical data  
31 into continuous data. Embedding tables can be as large as 1GB for NLP models, and embedding table  
32 look-ups are commonly implemented as row selection operations. However, because embedding table  
33 indices are from secret input data, embedded table accesses must be oblivious to input values in private  
34 inference. To implement embedding tables securely in the MPC setting, one can “densify” embedding  
35 table lookups by turning them into matrix multiplications. Unfortunately, our characterization study

36 shows that replacing an embedding table look-up with a matrix multiplication significantly increases  
 37 the execution time of embedding table operations.

38 For more efficient embedding table operations in MPC, we propose to utilize embedding tables’  
 39 compressibility. In particular, we use tensor train (TT) decomposition. The TT decomposition has  
 40 been applied to plaintext machine learning models [15, 16] to reduce the memory footprint and  
 41 speed up training processes. We exploit TT decomposition to reduce the size of matrices that encode  
 42 embedding tables and reduce the overhead for embedding table accesses. The experiments show that  
 43 the TT decomposition can speed up embedding table accesses by **2x** with only a 1.19 drop in the  
 44 masked-language model perplexity score.

## 45 2 Transformer-based Model MPC Inference Characterization

46 **Secure Multi-Party Computation.** In the MPC setting, an MPC client wishes to have remote servers  
 47 (non-colluding MPC servers) perform a sensitive task, such as translating classified sentences, without  
 48 revealing secret data and/or model parameters. The secret data is encrypted into two secret shares,  
 49 in a way that each share does not leak any information about the secret. There are two main secret  
 50 sharing formats that we use: additive and binary format. The binary secret sharing format is suitable  
 51 for bit-wise operations, whereas additive shares are more suitable for addition and multiplication.  
 52 Details about those formats can be found in [17]. After MPC servers receive their secret share, each  
 53 cloud server only sees and manipulates its own secret share, and returns the results to the client. The  
 54 client can combine the results from both servers to obtain the plaintext result. While MPC requires  
 55 that multiple parties do not collude and incur non-trivial communication overhead, MPC represents  
 56 one of the most promising techniques for PPML as its overhead is often much lower compared to HE.

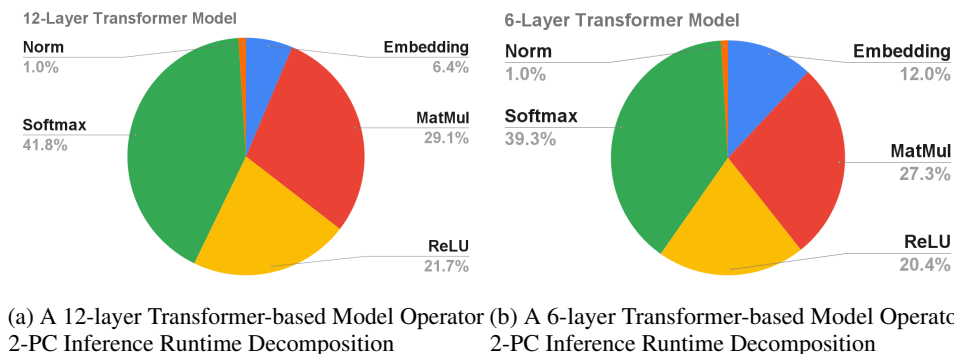


Figure 1: Transformer-based Model Operator 2-PC Inference Runtime Decomposition

57 **Experimental Setup.** In our characterization study, we obtained a secure 2-PC inference runtime.  
 58 We implemented MPC models using the CrypTen MPC framework [17]. The token embedding table  
 59 size is  $[250002 \times 1024]$ , and there are 6/12 layers of Transformers in the model ( $L = 6/12$ ,  $H =$   
 60  $1024$ ,  $A = 16$ ). Embedding tables are implemented as “densified” matrix multiplications. Runtimes  
 61 are obtained from using two nodes on the same server rack, and each node has an NVIDIA Tesla  
 62 V100 Volta GPU.

63 **MPC vs. Plaintext Inference.** Comparing with the plaintext single GPU inference, total inference  
 64 runtime is 12x slower with MPC. Embedding table accesses are 2, 523x slower; Matrix multiplications  
 65 are 3.5x slower; Softmax functions are 465x slower; and ReLU functions are 1, 372x slower. Among  
 66 all the operators, embedding tables and softmax sees the most overhead when MPC protocols are  
 67 implemented.

68 **MPC Execution Time Breakdown.** Figure 1a and Figure 1b show the breakdown of inference  
 69 runtime. Embedding table accesses make up 12% of total inference runtime. Softmax functions and  
 70 activation functions make up 40% and 28% of total inference runtime, respectively. Lightweight  
 71 non-linear functions such as ReLUs and Softmax dominate the inference runtime when MPC is used  
 72 because they need to be approximated with arithmetic or binary functions. Note that in the above  
 73 experiments the length of each input sentence is 128. If we increase the number of tokens in each  
 74 sentence, MPC softmax runtime grows more rapidly than all other operators because the input size

75 of softmax increases quadratically with sentence length. When token length of input sentences is  
 76 1024, softmax can account for 85% of the total inference runtime, and the total 2-party Transformer  
 77 inference is 316x slower than the computation in plaintext.

78 **Analysis of Softmax Overhead.** Among all the operations, softmax shows the largest slowdown  
 79 when MPC is applied. This slowdown is mainly due to a maximum function used in softmax for  
 80 numerical stability. Softmax functions for the  $i^{th}$  element in a vector size of  $n$  is defined as

$$Softmax(x_i) = \frac{e^{x_i}}{\sum_{k=1}^n e^{x_k}} \quad (1)$$

81 The exponential function is approximated using limit approximation [17]. The exponential function  
 82 can explode quickly even in plaintext, causing overflow when some input values are big. To achieve  
 83 numerical stability, softmax is practically implemented as

$$Softmax(x_i) = \frac{e^{x_i - x_{max}}}{\sum_{k=1}^n e^{x_k - x_{max}}} \quad (2)$$

84 where  $x_{max}$  is the maximum value in the given vector. The subtraction of the maximum value does  
 85 not change the final value of the softmax function, but it greatly improves the numerical stability. The  
 86 “maximum” function is typically cheap. However, in the context of MPC, the maximum function turns  
 87 into an extremely expensive operation. When using MPC, inputs are additively shared among multiple  
 88 parties. The maximum function requires a large number of comparisons, which are expensive in MPC  
 89 because comparisons require two secret sharing format conversions. This communication during  
 90 secret sharing format conversions makes maximum functions expensive, and there are  $\mathcal{O}(\log(N))$   
 91 comparisons. Our experiments show that if the maximum function is removed (for timing purpose  
 92 only) from the softmax operation, the softmax function becomes 7x and 9x faster on CPU and GPU,  
 93 respectively.

### 94 3 Tensor-Train Decomposition for Efficient Embedding Tables

95 **Tensor-Train (TT) Decomposition.** The basic idea of TT decomposition is to represent a big matrix  
 96 using tensor products of several smaller matrices. A tensor product is a function:

$$\mathbb{R}^{M_1 \times M_2 \dots \times M_k} \otimes \mathbb{R}^{N_1 \times N_2 \dots \times N_k} \rightarrow \mathbb{R}^{M_1 \cdot N_1 \times M_2 \cdot N_2 \dots \times M_k \cdot N_k} \quad (3)$$

97 Generalizing TT decomposition to an embedding matrix  $W \in \mathbb{R}^{M \times N}$ ,  $W$  can be decomposed into  $d$   
 98 smaller matrices  $w_k \in \mathbb{R}^{R_{k-1} \times m_k \times n_k \times R_k}$ , where  $M = \prod_{k=1}^d m_k$ ,  $N = \prod_{k=1}^d n_k$ , and  $R_0 = R_d =$   
 99 1. We refer  $R_k$  as the ranks of decomposed matrices. For example, if  $d$  is 2, an embedding size  
 100 of  $[250002 \times 1024]$  can be decomposed into two smaller matrices of  $[1 \times 500 \times 32 \times rank]$  and  
 101  $[rank \times 502 \times 32 \times 1]$ . Note that  $502 \times 500 > 250002$  and  $32 \times 32 = 1024$ . If the rank of the all  
 102 decomposed matrices are 4, the original matrix with 24M parameters can be decomposed into two  
 103 smaller 64K-element matrices. When accessing certain locations in the original embedding table, a  
 104 entry in every decomposed matrices is fetched (this fetching is implemented as dense one-hot matrix  
 105 multiplications). To reconstruct the original entry, apply dot products at the dimension of the rank  
 106 among fetched entries.

107 **TT Decomposition for Embedding Tables in MPC.** TT decomposition can be applied to the matrix  
 108 multiplications that are used for embedded tables to reduce the overhead, enabling a trade-off  
 109 between performance and model accuracy. Here, we present the embedding table query runtime and  
 110 model accuracy using different configurations of TT decomposition. We use  $d/ranks$  to represent  
 111 TT decomposition configurations.  $d$  represents the number of smaller decomposed matrices, and  
 112  $ranks$  represents the rank of each decomposed matrix. Configuration 3/64 means that the original  
 113 embedding table is decomposed into 3 smaller 64-rank matrices. Throughout our experiments, the  
 114 original embedding table size is  $[250002 \times 1024]$ . With fixed configurations, the dimensions of  
 115 smaller matrices using different configurations of TT decomposition are shown in Table 1. The  
 116 performance results are obtained by running 2-party inference using CrypTen [17] on two nodes on  
 117 the same server rack, each node with an NVIDIA Tesla V100 Volta GPU. To evaluate the impact  
 118 of model accuracy, the experiments on a masked language model are performed without MPC on a  
 119 single node with 8 NVIDIA Tesla V100 Volta GPUs.

$3/ranks$	$4/ranks$	$5/ranks$
$1 \times 50 \times 8 \times ranks$	$1 \times 20 \times 4 \times ranks$	$1 \times 10 \times 4 \times ranks$
$ranks \times 65 \times 8 \times ranks$	$ranks \times 24 \times 4 \times ranks$	$ranks \times 10 \times 4 \times ranks$
$ranks \times 80 \times 16 \times 1$	$ranks \times 25 \times 8 \times ranks$	$ranks \times 10 \times 4 \times ranks$
-	$ranks \times 25 \times 8 \times 1$	$ranks \times 13 \times 4 \times ranks$
-	-	$ranks \times 20 \times 4 \times 1$

Table 1: Matrix Dimensions using TT Decomposition

Number of Matrices	3			4			5		
Ranks	64	128	196	64	128	196	64	128	196
Batch=32 Speedup	16.37	6.95	2.91	16.68	5.68	3.35	17.88	6.11	2.91
Batch=64 Speedup	11.93	3.84	2.09	9.28	3.15	1.34	8.01	3.40	1.44
Batch=128 Speedup	8.65	2.90	1.18	6.34	1.99	0.88	6.98	1.78	0.88
Batch=256 Speedup	6.56	1.79	0.90	4.67	1.41	0.64	3.97	1.39	0.65
Batch=512 Speedup	4.95	1.19	0.56	2.98	0.93	0.40	2.50	0.95	0.40

Table 2: TT Decomposition Inference Runtime Improvement

120 Table 2 demonstrates embedding table inference speedups using different TT decomposition configu-  
121 rations with batch sizes. The batch size represents the number of embedded table accesses that are  
122 performed together. When batched accesses are small, most configurations demonstrate speedups.  
123 However, with more batched accesses, the configurations with more decomposed matrices and more  
124 ranks begin to show negative speedups. Computations needed and the number of bytes of data  
125 communicated for reconstructing embedding entries from decomposed matrices increases linearly  
126 with the numbers of accesses. On the contrary, the baseline’s one-hot matrix multiplication’s runtime  
127 does not grow linearly because the number of communication rounds is amortized as the number of  
128 batched accesses grows (the size of the embedding table does not depend on batch sizes). The results  
129 suggest that TT decomposition can significantly improve performance for cases with small batch  
130 sizes as in real-time inference settings. For example, users’ inputs to real-time translation software  
131 are generally no longer than 50 tokens where one token corresponds to one embedding table access;  
132 length of 90% data points in MNLI[18], XNLI[19] and CoNLL-2003 [20] is smaller than 75.

133 The speedups and compression come with a cost. Besides inference runtime reduction, we also  
134 measured the TT decomposition’s impact on model accuracy. We have run a masked language model  
135 on WikiText-103 [21]. The Transformer configuration we used is ( $L = 24$ ,  $H = 1024$ ,  $A = 16$ ),  
136 and the token embedding table size is  $[250002 \times 1024]$ . We have trained the model for ten epochs  
137 and reported their best perplexity score. Table 3 presents perplexity of various TT decomposition  
138 configurations. Configuration 3/64 achieves a speedup of 2.09x, while incurring a 1.19 loss in  
139 perplexity score. If applications can tolerate higher loss in perplexity scores, using TT decomposition  
140 configurations such as (3/64, 3/128, 4/64) can achieve even more speedups.

## 141 4 Conclusions

142 Transformers are widely used in NLP tasks and will likely represent an important workload for PPML  
143 in the future. Interestingly, this paper shows that Transformers introduce new research challenges that  
144 do not exist for private inference of CNNs. While this paper shows that TT decomposition can be  
145 used to speed up embedding table accesses, further optimizations, especially for non-linear operations  
146 such as softmax, are needed to enable private real-time NLP inference.

Configs	Plaintext	3/64	3/128	3/196	4/64	4/128	3/196	5/64	5/128	5/196
PPL	12.8	17.0	15.17	14.04	18.25	16.47	14.41	17.82	16.09	14.83

Table 3: Masked Language Model Perplexity Score

## References

- 147
- 148 [1] Nishant Kumar, Mayank Rathee, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul  
149 Sharma. Cryptflow: Secure tensorflow inference. In *IEEE Symposium on Security and Privacy*.  
150 IEEE, May 2020.
- 151 [2] Sameer Wagh, Shruti Tople, Fabrice Benhamouda, Eyal Kushilevitz, Prateek Mittal, and Tal  
152 Rabin. Falcon: Honest-majority maliciously secure framework for private deep learning.  
153 *Proceedings on Privacy Enhancing Technologies*, 2020.
- 154 [3] Minsu Cho, Zahra Ghodsi, Brandon Reagen, Siddharth Garg, and Chinmay Hegde. Sphynx:  
155 Relu-efficient network design for private inference. 2021.
- 156 [4] Nandan Kumar Jha, Zahra Ghodsi, Siddharth Garg, and Brandon Reagen. Deepreduce: Relu  
157 reduction for fast private inference. *Proceedings of the 38 th International Conference on*  
158 *Machine Learning*, 2021.
- 159 [5] Zahra Ghodsi, Nandan Kumar Jha, Brandon Reagen, and Siddharth Garg. Circa: Stochastic  
160 relus for private deep learning. *arXiv preprint arXiv:2106.08475*, 2021.
- 161 [6] Xiaoqiang Sun, Peng Zhang, Joseph K. Liu, Jianping Yu, and Weixin Xie. Private machine  
162 learning classification based on fully homomorphic encryption. *IEEE Transactions on Emerging*  
163 *Topics in Computing*, 8(2):352–364, 2020.
- 164 [7] Simon Johnson, Vinnie Scarlata, Carlos Rozas, Ernie Brickell, and Frank Mckeen. Intel software  
165 guard extensions: Epid provisioning and attestation services. 2016.
- 166 [8] David Lie Chandramohan Thekkath, Mark Mitchell, Patrick Lincoln, Dan Boneh, John Mitchell,  
167 and Mark Horowitz. Architectural support for copy and tamper resistant software. *ACM*  
168 *SIGARCH Computer Architecture News*, 2000.
- 169 [9] ARM Limited. Arm security technology building a secure system using trustzone technology.  
170 2016.
- 171 [10] Florian Tramèr and Dan Boneh. Slalom: Fast, verifiable and private execution of neural  
172 networks in trusted hardware. *arXiv preprint arXiv:1806.03287*, 2019.
- 173 [11] Krishna Giri Narra, Zhifeng Lin, Yongqin Wang, Keshav Balasubramaniam, and Murali An-  
174 navaram. Privacy-preserving inference in machine learning services using trusted execution  
175 environments. *IEEE International Conference on Cloud Computing*, 2021.
- 176 [12] Hanieh Hashemi, Yongqin Wang, and Murali Annavaram. Darknight: A data privacy scheme  
177 for training and inference of deep neural networks. *Proceedings on the 54th International*  
178 *Symposium on Microarchitecture*, 2021.
- 179 [13] Oded Goldreich. Secure multi-party computation. 1998.
- 180 [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
181 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly  
182 and Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for  
183 image recognition at scale. *Proceedings of Ninth ICLR*, 2021.
- 184 [15] Oleksii Hrinchuk, Valentin Khrulkov, Leyla Mirvakhabova, Elena Orlova, and Ivan Oseledets.  
185 Tensorized embedding layers for efficient model compression. *Proceedings of Ninth ICLR*,  
186 2021.
- 187 [16] Chunxing Yin, Bilge Acun, Xing Liu, and Carole-Jean Wu. Tt-rec: Tensor train compression  
188 for deep learning recommendation model embeddings. *arXiv preprint arXiv:2101.11714*, 2021.
- 189 [17] B. Knott, S. Venkataraman, A.Y. Hannun, S. Sengupta, M. Ibrahim, and L.J.P. van der Maaten.  
190 Crypten: Secure multi-party computation meets machine learning. In *arXiv 2109.00984*, 2021.
- 191 [18] Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus  
192 for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2018.

- 193 [19] Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R. Bowman, Holger  
194 Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. *arXiv*  
195 *preprint arXiv:1809.05053*, 2018.
- 196 [20] Erik F., Tjong Kim Sang, and Fien De Meulder. Introduction to the conll-2003 shared task:  
197 Language-independent named entity recognition. *arXiv preprint arXiv:0306050*, 2003.
- 198 [21] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
199 models. *arXiv preprint arXiv:1609.07843*, 2016.