

SecNDP: Secure Near-Data Processing with Untrusted Memory

Wenjie Xiong*[§], Liu Ke*^{†§}, Dimitrije Jankov[‡], Michael Kounavis*, Xiaochen Wang*, Eric Northup*, Jie Amy Yang*, Bilge Acun*, Carole-Jean Wu*, Ping Tak Peter Tang*, G. Edward Suh*[◊], Xuan Zhang[†], Hsien-Hsin S. Lee*

*Meta, [†]Washington University in St. Louis, [‡]Rice University, [◊]Cornell University
 wenjiex@fb.com, ke.l@wustl.edu, dj16@rice.edu,

{michaelkounavis, xiaochenwang, digitaleric, amyyang, acun, carolejeanwu, ptpt, edsuh}@fb.com,
 xuan.zhang@wustl.edu, lehs@fb.com

Abstract—Today’s data-intensive applications increasingly suffer from significant performance bottlenecks due to the limited memory bandwidth of the classical von Neumann architecture. Near-Data Processing (NDP) has been proposed to perform computation near memory or data storage to reduce data movement for improving performance and energy consumption. However, the untrusted NDP processing units (PUs) bring in new threats to workloads that are private and sensitive, such as private database queries and private machine learning inferences. Meanwhile, most existing secure hardware designs do not consider off-chip components trustworthy. Once data leaving the processor, they must be protected, *e.g.*, via block cipher encryption. Unfortunately, current encryption schemes do not support computation over encrypted data stored in memory or storage, hindering the adoption of NDP techniques for sensitive workloads.

In this paper, we propose *SecNDP*, a lightweight encryption and verification scheme for untrusted NDP devices to perform computation over ciphertext and verify the correctness of linear operations. Our encryption scheme leverages arithmetic secret sharing in secure Multi-Party Computation (MPC) to support operations over ciphertext, and uses counter-mode encryption to reduce the decryption latency. The security of the encryption and verification algorithm is formally proven. Compared with a non-NDP baseline, secure computation with SecNDP significantly reduces the memory bandwidth usage while providing security guarantees. We evaluate SecNDP for two workloads of distinct memory access patterns. In the setting of eight NDP units, we show a speedup up to $7.46\times$ and energy savings of 18% over an unprotected non-NDP baseline, approaching the performance gain attained by native NDP without protection. Furthermore, SecNDP does not require any security assumption on NDP to hold, thus, using the same threat model as existing secure processors. SecNDP can be implemented without changing the NDP protocols and their inherent hardware design.

Keywords—Security and Privacy, Near-Data Processing, Cryptography, Privacy-Preserving Machine Learning

I. INTRODUCTION

The classical von Neumann architecture separates computation from data storage. However, memory and interconnect bandwidth have not been able to keep up with the scaling out and scaling up of the processor cores, hitting the (in)famous *memory wall* [77]. For data-intensive

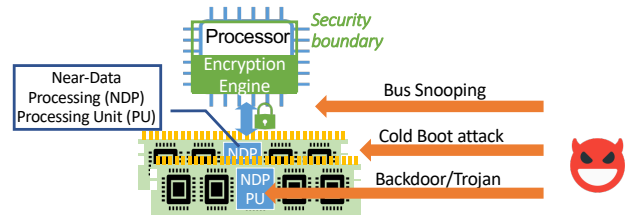


Figure 1. Threat model of SecNDP.

applications with high parallelism and localized memory accesses, memory bandwidth has become a major performance bottleneck. Therefore, techniques such as near-data processing (NDP) offer much promise in both performance and energy consumption by offloading computation to main memory (*i.e.*, DRAM) [8], [10], [24], [36], [44] or even storage [45], [64], [76]. In Samsung’s Aquabolt-XL (HBM2-PIM), for example, addition and multiplication operators are supported in DRAM [46].

At the same time, many data-intensive applications work on sensitive data, such as queries on private databases (*e.g.*, medical history, user information) and machine learning inference using private models (*e.g.*, models that need IP protection [17] or may reveal the private training dataset [68]). As more computation is outsourced to the cloud, sensitive workloads are becoming more susceptible to cyber and physical attacks. In addition to the pursuit of better performance and lower energy, guaranteeing security and privacy is now considered a first-class citizen in designing future computing systems. To protect the confidentiality and/or integrity of program execution, secure hardware systems designed with a Trusted Execution Environment (TEE) [1], [2], [4], [14], [73] have been proposed and implemented in commercial processors to enforce protection. All of these secure hardware designs consider off-chip components to be untrusted (in Figure 1) – be it a TEE design in CPU [14], [73], CPU and GPU [75], or other accelerators [11], [33]. In this threat model, off-chip memory needs to be encrypted for confidentiality and a message authentication code (MAC) for each data block in memory needs to be computed and checked for integrity.

However, such existing memory protection mechanisms prevent the adoption of NDP in a TEE, because the NDP

[§] The first two authors contributed equally to this work.

Processing Units (PUs) only have access to encrypted data, a.k.a., ciphertext. One may apply homomorphic encryption (HE) [23], [59] to enable computations over ciphertext. However, the state-of-the-art HE incurs at least four orders of magnitude performance slowdown [60]. The overhead of HE is too high to let NDP outperform a TEE without NDP. Another approach for secure NDP is to include the NDP PUs inside the TEE [5], [9]. However, this approach requires trusting multiple hardware vendors, an extra key exchange protocol, and hardware components on the NDP PUs, resulting in an undesirably large attack surface for the trusted computing base (TCB). There is no existing practical solution to achieve secure NDP for untrusted memory.

In this paper, we propose SecNDP – a lightweight encryption and verification scheme for a secure processor to use untrusted NDP. SecNDP uses a secure Multi-Party Computation (MPC) protocol [20] between the processor and the untrusted memory, where a block cipher generates the processor’s share of secret on-chip without additional off-chip accesses. Like MPC, SecNDP encryption supports addition and scale multiplication efficiently. Like counter-mode encryption [72], SecNDP has low decryption latency with only one adder on the performance critical path. Furthermore, SecNDP introduces a verification tag based on a linear checksum [42] to verify the correctness of the results from the untrusted NDP PU. To efficiently implement SecNDP, we propose architectural support in a secure processor. The SecNDP design only requires relatively small changes in the processor while leaving the NDP protocols and hardware intact. SecNDP’s performance and accuracy are evaluated through detailed cycle-level simulations using two real-world data-intensive workloads, *deep learning recommendation inference* and *medical data analytics*. The results show that SecNDP is up to $7.46\times$ faster while using 18% less memory energy compared to the unprotected non-NDP baseline, approaching the performance and energy efficiency of native NDP without protection. The following summarizes the key contributions:

- We propose the first encryption and verification schemes that enable practical and provably secure computation in untrusted NDP.
- We tailor-design the SecNDP architecture, to demonstrate low decryption and verification latency, and low area and power overhead on the processor. The design does not modify NDP PUs and protocols.
- We demonstrate SecNDP’s performance benefits and accuracy on two real-world data-intensive use cases.

II. THREAT MODEL

SecNDP is to be used with a TEE design. We follow the threat model of a typical TEE, shown in Figure 1. The processor (CPU, GPU, or an accelerator) is trusted and considered secure, *i.e.*, there is no hardware backdoor. An attacker cannot observe or manipulate the internal state of

the processor. We assume an attacker’s software co-located in the processor cannot access protected data (directly or through side channels) within the processor. A privileged attacker may access unencrypted data in memory. Data buses are susceptible to passive eavesdropping and unauthorized modifications. For volatile memory such as DRAM, the attacker can conduct a cold-boot attack [29] to dump the memory content. For non-volatile memory, the attacker can unplug the storage to access and modify the content. The NDP processing units (PUs) are untrusted. NDP PUs may have a backdoor or a Trojan to leak data or return a malicious computation result.

III. BACKGROUND

A. Near-Data Processing (NDP)

For data-intensive applications, such as machine learning [10], [36], [44], [51], [76], similarity search [47], data base [35], graph processing [6], and stencil computing [70], memory bandwidth has become the major performance bottleneck. Many of such workloads deal with vectors and matrices and conduct linear operations, which is our focus.

Depending on the data reuse rate, moving data from memory or storage to CPU for computation is inefficient [69]. Many prior works propose near-memory processing using 3D/2.5D-stacked DRAM technology (such as HMC/HBM) [7], [38], [43], [57]. Accelerators inside the data buffer devices in a commodity DIMM are also proposed to support large-scale workloads [8]. For specific workloads such as sparse embedding operations in recommendation models, tailor-designed light-weight NDP systems are proposed to leverage rank-level parallelism with higher speedup potential [10], [36], [37], [44]. In addition to near-memory processing, near-storage processing is proposed to process data near storage with larger capacity [45], [64], [76]. With the accelerated development of 3D stacking and emerging memory technologies compatible with CMOS process [49], [52], productizing NDP on a computing platform is closer to reality than ever [7], [37], [43], [46].

B. TEEs and Memory Protection

In hardware-based TEE designs, such as AEGIS [73], Mi6 [14], as well as commercial solutions including Intel SGX [4], AMD SEV [1], and ARM TrustZone [2], off-chip memory is untrusted. Off-chip data must be encrypted for confidentiality protection and data integrity needs to be verified.

Confidentiality Protection with Counter-Mode Encryption. In counter-mode encryption [42], [54], [72], a data block (e.g., a cache line) is eXclusive-ORed (XOR) with a One-Time Pad (OTP) (Figure 2(a)). The OTP is an encrypted counter x_e generated using a block cipher such as Advanced Encryption Standard (AES) with the address of the data block and a unique version number v as inputs. The cipher uses the processor’s secret key K . The security

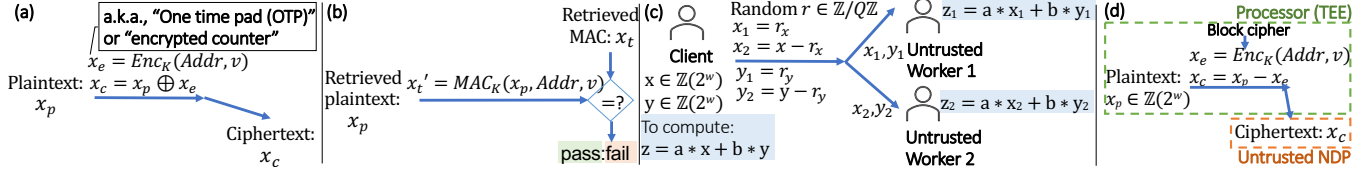


Figure 2. (a) Counter mode encryption. (b) Integrity verification using MAC. (c) Arithmetic secret sharing in secure MPC. (d) Proposed SecNDP.

of the counter-mode encryption relies on the uniqueness of the version number v for each address. If the version number is reused for multiple plaintext blocks at the same address, the attacker can learn the relationship between the plaintext blocks.

When the processor fetches a data block from memory, the OTP can be computed using the address and the version number in parallel with the data access. When the ciphertext is returned, only one XOR with the OTP is needed for decryption. Counter-mode encryption is widely used for memory protection due to its functional parallelism resulting in lower latency [67], [72].

Integrity using Message Authentication Codes (MACs).

To detect an unauthorized modification, a keyed message signature, a.k.a., Message Authentication Code (MAC), is stored in memory for each data block. When the processor fetches a data block from memory, its MAC is also fetched. Independently, we use the secret key to compute the MAC to see if it matches with the retrieved one (Figure 2(b)). Without the secret key, the probability that an attacker generates the corresponding MAC for an altered message is negligible. Thus, any change in data from memory will result in a MAC mismatch. To prevent attacks from replaying stale data values or copying valid data from a different address with valid MAC, the address $Addr$ and version number v are also incorporated into the MAC. The integrity verification passes only when both the $Addr$ and v also match. To protect the integrity of the version numbers v , these values are kept on-chip or protected with an integrity tree [62]. To protect both confidentiality and integrity, authenticated encryption with associated data (AEAD) schemes such as the AES GCM mode [54] and the AES CWC mode [42] are designed to use one secret key to both encrypt data and generate a MAC.

Linear Checksum and MACs. A hash function or checksum maps messages into short outputs. A universal hash has the property that the probability of any given pair of messages having the same checksum is small. Such a universal hash function can be used as the building block for a MAC. Linear Modular Hashing [30] is an almost universal hash function suitable for a fast software or hardware implementation. It is applicable to variable-sized data and used in fast message authentication, such as the AES CWC mode [42]. In this paper, we use a linear modular hash not only for its performance but also to leverage its linearity to verify linear operations in NDP.

C. Secure Multi-Party Computation (MPC)

One popular approach to outsource sensitive workloads to untrusted workers is to use secret sharing [20], [55]. MPC secret sharing splits a secret into multiple shares and distributes them to workers. Each worker performs computation over its share locally. Assuming the workers will not collude, it is information-theoretically impossible for each worker to recover the secret from its share. SecNDP let the TEE hold one share and the memory hold another share, and thus, do not require the non-collusion assumption to hold among NDP PUs.

Arithmetic secret sharing [20], [55] is one of the most widely used secret sharing schemes. Let $\mathbb{Z}(2^{w_e})$ denote the integer ring of size 2^{w_e} . The shares are constructed such that the sum of all shares is equal to the original secret value $x \in \mathbb{Z}(2^{w_e})$. Figure 2(c) shows arithmetic sharing between two workers. Each worker holds a share of x , where shares are denoted by $x_1 \in \mathbb{Z}(2^{w_e})$ and $x_2 \in \mathbb{Z}(2^{w_e})$, and $x = x_1 + x_2$. Arithmetic secret sharing enables ciphertext-to-ciphertext additions and additions/multiplications by constants over ciphertext. For example, to compute $z = a * x + b * y$, each worker performs an addition on their respective shares of x and y . In the end, the client collects the two partial results z_1 and z_2 , and adds them together to obtain the final result. Compared to HE schemes on high-dimensional rings [60], the operations for each worker are simple and lightweight. This is with arithmetic secret sharing.

All the operations of arithmetic sharing are in the integer ring $\mathbb{Z}(2^{w_e})$. Thus, if the original data are in a floating-point format, they must be quantized into fixed-point numbers or integers. For machine learning models, previous studies have shown that quantization does not have significant impact on the accuracy of certain ML models [31], [41], [74]. More recently, MPC in floating-point arithmetic has also been explored [27], but such evaluation is outside the scope of this work.

IV. SECNDP ENCRYPTION AND VERIFICATION SCHEMES

In this section, we describe an arithmetic encryption and verification scheme. The encryption supports computation over ciphertexts while providing confidentiality guarantees. The verification scheme further protects the integrity of data using extra data tags. The encryption scheme can be used alone without verification, or with verification for stronger security.

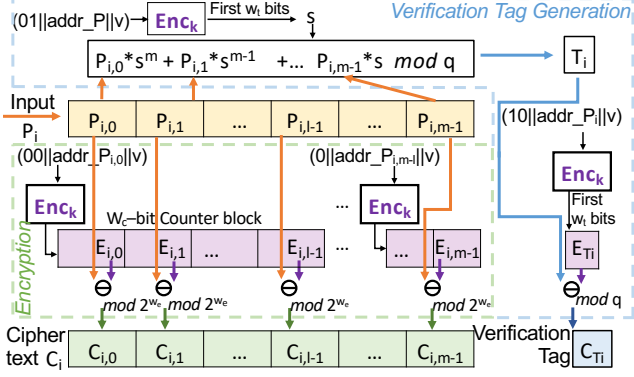


Figure 3. Diagram of the proposed arithmetic encryption and verification tag generation scheme. Here, l is the number of data elements in a data block of the block cipher size (i.e., $l = w_c/w_e$). \ominus represents subtraction in the ring $\mathbb{Z}(2^{w_e})$. $E_{i,j}$ is the j -th OTP substring associated with row i .

A. Preliminaries and Notations

We use vector-matrix multiplication as an example operation as it is the most critical and widely-used kernel in machine learning and other data-intensive applications. Section VI-A introduces two use cases that use variants of vector-matrix multiplication. Our scheme can be generalized for other operations as well. We use P to denote the plaintext of a 2-D matrix of size $n \times m$, $P_{i,j}$ to denote a element in the array, p or P_i to denote a row vector in P , and a to denote a vector of dimension n . Each element in p and a is an integer (or a fixed point number) of width w_e . We assume w_e to be a power of 2 and smaller than a cache line size. We use $E(K, X)$ to denote the output of a w_c -bit block cipher (e.g., AES) that encrypts a w_c -bit message X using a w_K -bit secret key K . The elements of the matrix are also considered as consisting of chunks that are w_c bits long. Each chunk contains $l \leftarrow w_c/w_e$ elements.

B. SecNDP Arithmetic Encryption

SecNDP arithmetic encryption is a combination of arithmetic secret sharing and counter-mode encryption, as shown in Figure 2 (d). The scheme details are depicted in Algorithm 1 and Figure 3. First, the plaintext P will be divided into cnt_c binary strings (chunks) of size w_c . The starting physical address of each chunk and the version number v will be used as the input to the block cipher, which generates an OTP using a standard block cipher (e.g., AES). The concatenation of all cnt_c OTP blocks (e_0, e_1, \dots) will have the same number of bits as the plaintext P .

Then, the plaintext P and the concatenation of all the OTP blocks e are divided into binary strings of size w_e . The j -th strings are denoted by p_j and e_j , respectively. The ciphertext c_j is the plaintext p_j subtracted by the corresponding OTP block e_j in the ring $\mathbb{Z}(2^{w_e})$. We name the scheme “Arithmetic Encryption Scheme”, because the

Algorithm 1: Arith. Encryption, $Arith-E(K, P, Addr)$

```

1 Inputs:  $P, K, Addr$ ; //  $Addr$  is the address of  $P$ 
2  $v \leftarrow \mathcal{V}()$ ; // unique version, also considered padded with zeros
3  $cnt_c \leftarrow size(P)/w_c$ ; // number of blocks of block cipher size
4 for  $i = 0$  to  $cnt_c - 1$  // each data block
5 do
6    $Addr_i \leftarrow Addr + i \times (w_c/8)$ ; // addr. of block  $i$  in bytes
7    $e_{Addr_i} \leftarrow E(K, 00 || Addr_i || v)$ ; // OTP of block  $i$ 
8   for  $j = 0$  to  $(w_c/w_e) - 1$  // for every  $w_e$ -bit element
9   do
10     $e_j \leftarrow e_{Addr_i}[j \times w_e : (j+1) \times w_e]$ ; // OTP of element  $j$ 
11     $p_j \leftarrow P[Addr_i + j \times w_e : Addr_i + (j+1) \times w_e]$ ;
12     $c_j \leftarrow p_j - e_j \pmod{2^{w_e}}$ ; // ciphertext of element  $j$ 
13  end
14   $C_{Addr_i} \leftarrow c_0 || \dots || c_{w_c/w_e - 1}$ ; // concatenating ciphertexts
15 end
16 Return concatenation of all  $C_{Addr_i}$ ;

```

ciphertext c_j and the OTP block e_j can be seen as the arithmetic shares of the secret p_j .

C. Computation over Ciphertext in SecNDP

After the arithmetic encryption completes, each of the processor and the NDP device holds an arithmetic share of the secret. They can then follow the standard MPC protocol [20], [55] to conduct computation. With arithmetic secret sharing, additive operations such as addition and multiplication with a non-private scale are lightweight. In this paper, we demonstrate how to perform a multiplication with a non-private vector.

Figure 4 (a) shows an example of a non-private vector a multiplied with a private matrix P . In the initialization step (T0), the matrix P is encrypted by Alg. 1 and the resulting ciphertext C is stored in untrusted memory or data storage. When there is a vector a to compute $a \times P$ (T1 in Figure 4), the untrusted NDP will compute $a \times C$. The NDP operation is exactly the same as the one in the unprotected scenario. Meanwhile, the processor will use the address and the version number v to generate the OTP blocks E for the entire plaintext P and compute $a \times E$. In the end, the NDP PU sends its share of result back to the processor, and the processor adds the two shares for the final result.¹ The encryption scheme has the property that $E + C = P$, and thus $res = C_{res} + E_{res} = a \times C + a \times E = a \times P$.

D. Comparing SecNDP with TEE and Unprotected NDP

Vector-matrix multiplication in a TEE and unprotected NDP are also shown in Figure 4. All of them have the initialization step (T0) that initializes data and loads data to memory. In the TEE and SecNDP, the data is encrypted before being loaded to memory. Regarding the computations performed, when compared to a TEE baseline, SecNDP

¹Note that all the operations are in the ring $\mathbb{Z}(2^{w_e})$, and thus, overflow could happen. However, overflow can be detected by our verification scheme.

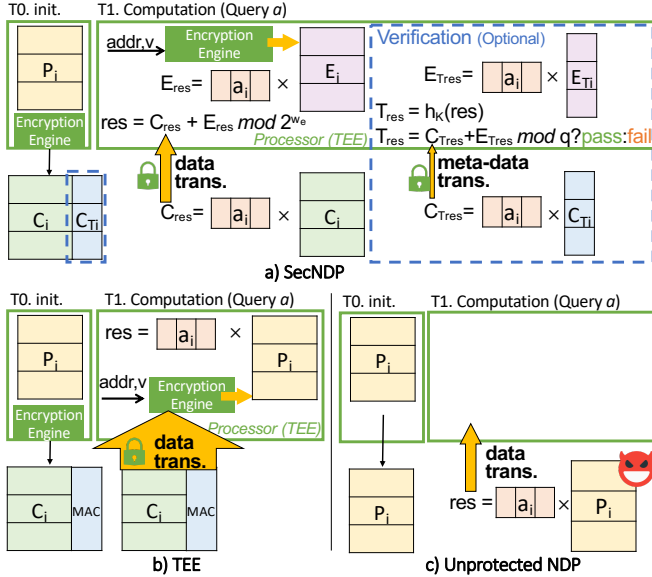


Figure 4. Vector-matrix multiplication in SecNDP, TEE, and unprotected NDP. NDP and SecNDP significantly save memory bandwidth usage, which is the bottleneck for data-intensive workloads. Green boxes indicate encrypted or protected content. Blue boxes show the verification tags.

only needs to transfer the results back to the processor, and thus reduces memory bandwidth usage and saves energy. When there are multiple NDP PUs, they can access memory and perform the computations in parallel to improve performance. However, since both the processor and NDP need to compute on their own share of secret, SecNDP does not save computation on the processor side. Nonetheless, for data-intensive workloads, memory bandwidth, rather than computation, is the main bottleneck. Hence, performance can be significantly improved by SecNDP over the non-NDP baseline.

Compared to an unprotected NDP, the off-chip operations and the required data movement are exactly the same. Hence, there is no modification in the NDP implementation needed to support SecNDP. For decryption and verification, the OTP (E_{res}) and verification tag ($E_{T_{res}}$) computations can be done in parallel with the NDP operations. Only one extra final addition is added to the performance critical path (more details in Section V-E3), which results in negligible performance overhead.

E. Security of the Arithmetic Encryption

The arithmetic encryption applies two-party arithmetic secret sharing between the TEE and the memory, with the TEE's share of secret (*i.e.*, OTP) generated from a block cipher. If an attacker can distinguish c_j (in Alg. 1) from a random number, an oracle can be built to break the block cipher. The security is similar to that of counter-mode encryption [42], [54]².

²We have the proof and formal definitions in the appendix of the full paper at <https://eprint.iacr.org/2021/1642>.

Let $E_D(K, Addr, v)$ denote the randomized encryption systems $E(K, D || Addr || v)$. The secret key $K \in \{0, 1\}^{w_K}$ is drawn from the uniform distribution. D is a binary string, which can be one of '01', '00' or '10' (in Alg. 1, 2, and 3, respectively). We consider that version values $v \in \{0, 1\}^{w_v}$ are drawn from distribution $\mathcal{V}()$ and padded with zeros. Version numbers are not in the control of any adversary but drawn by algorithms 1, 2 and 3. Let $\text{Adv}_{|Q|}^{E_D(\cdot)}$ be the distinguishing advantage³ associated with $E_D(\cdot)$ and query budget $|Q|$, where version values are unique per distinct input and not in the control of the distinguisher.

Theorem 1 (Arithmetic Encryption is Secure). *Let $\text{Arith-E}(K, P, Addr)$ be the arithmetic encryption system of Alg. 1. Let also $\mathcal{A}_{\text{CPA}}^{\text{Arith-E}(\cdot), Q_e}$ be an adaptive chosen plaintext adversary with query budget $|Q_e|$ attacking $\text{Arith-E}(K, P, Addr)$, where plaintext is P , $Addr$. Then, the advantage of this adversary is bounded in the following way:*

$$\text{Adv}(\mathcal{A}_{\text{CPA}}^{\text{Arith-E}(\cdot), Q_e}) \leq \frac{1}{2^{w_K}} + \text{Adv}_{|Q|'}^{E_{00}(\cdot)} \quad (1)$$

where $|Q|' = \frac{m \cdot n \cdot w_e}{w_c} \cdot |Q_e|$.

The version number v should be generated securely and have its integrity protected. We will discuss how v is managed in Section V-A. The security level depends on the key size w_K and the block cipher. Block ciphers such as AES are considered to be well designed pseudo-random permutations, and their encryption output is practically indistinguishable from random output. Thus, if $E_{00}(\cdot)$ is based on AES, $\text{Adv}_{|Q|'}^{E_{00}(\cdot)}$ is negligible.

F. SecNDP Verification

Since the NDP is not trusted, the secure processor also needs to verify the correctness of the computation (including overflow). Here, we propose a verification scheme for linear operations. We use a MAC-then-encrypt strategy [61] to construct a verification tag. In the case of linear operations, we choose Linear Modular Hashing [30], [42] as the MAC ($h_K(\cdot)$ in Alg. 2) and apply SecNDP arithmetic encryption to encrypt the resulting MAC (in Alg. 3 and Figure 3). Here, q is a big prime number, *e.g.*, $2^{127} - 1$ which is w_t bits long. The memory side stores a verification tag C_{T_i} for each row vector P_i in P .

To verify the computation results, the secure processor computes a MAC of the result using the secret key and matches that with the retrieved MAC, similar to memory integrity protection in Figure 2(b). Here, the NDP and the processor need to compute together to retrieve the MAC leveraging arithmetic encryption. Specifically, as shown in Figure 4 (a), the NDP computes the result of a vector multiplication over the encrypted checkcodes to get an encrypted

³Distinguishing advantage of block cipher $E_D(\cdot)$ is the probability that an attacker can distinguish the output of $E_D(\cdot)$ from a truly random output, *i.e.*, breaking the block cipher.

Algorithm 2: Linear Checksum, $h_K(P_i)$

1 **Inputs:** K, P_i , **Output:** T_i
2 $P \leftarrow$ matrix containing P_i ; $paddr(P) \leftarrow$ Address of P ;
3 $v \leftarrow \mathcal{V}()$; //version padded with zeros, drawn once for P
4 $s \leftarrow$ first w_t bits of $E(K, 01 || paddr(P) || v)$;
5 **Return** $T_i \leftarrow \sum_{j=0}^{m-1} P_{i,j} \times s^{(m-j)} \bmod q$;

Algorithm 3: Encrypted MAC, $el\text{-}MAC(K, P_i, Addr_i)$

1 **Inputs:** $K, P_i, Addr_i$ //Address of P_i , **Output:** C_{T_i}
2 $T_i \leftarrow h_K(P_i)$; //linear checksum, coming from Alg. 2.
3 $v \leftarrow \mathcal{V}()$; //version padded with zeros, drawn once for P
4 $E_{T_i} \leftarrow$ first w_t bits of $E(K, 10 || paddr(P_i) || v)$; //OTP for tag
5 **Return** $C_{T_i} \leftarrow T_i - E_{T_i} \bmod q$; // ciphertext of tag;

tag $C_{T_{res}} = a \times C_T$. The processor computes $E_{T_{res}} = a \times E_T$, and $C_{T_{res}} + E_{T_{res}}$ will be used as the retrieved MAC. Due to the linearity of h_K , $h_K(res) = h_K(a \times P) = a \times h_K(P) = a \times (C_T + E_T) = C_{T_{res}} + E_{T_{res}}$ holds.

As shown in Figure 4(a), our verification scheme associates a tag C_{T_i} with each row of the matrix. Like other memory integrity protection schemes, our scheme requires extra memory to store the verification tag. However, because the tag is for each row, the tag is relatively small compared to the data when the matrix has large row sizes. Furthermore, we use arithmetic encryption for tags. Thus, NDP computes and returns the tag of the result $C_{T_{res}} = a \times C_T$, and not all the tags of the rows that participate in the vector matrix multiplication operation, reducing the memory bandwidth requirement.

G. Security of the Verification

The security of the verification scheme is derived from the fact that the verification tags are always encrypted on the memory side and s remains to be a secret to the untrusted memory. In order to generate the valid verification tag for a new value, an attacker needs to correctly guess s . Replay attacks are prevented by including a version number in the MAC (Alg. 2).

Theorem 2 (Integrity of the Weighted Summation). *Let $ws\text{-}MAC_K(P, Addr)$ and $ws\text{-}Verify_K(C, Addr)$ be the MAC and verification oracles associated with the weighted summation operation (i.e., vector matrix multiplication). Let also $\mathcal{A}_{MAC}^{ws\text{-}MAC(), Q_s, Q_v}$ be a standard MAC adversary, playing a MAC forgery game on weighted summation oracles. Q_s is the set of sign queries issued by the MAC adversary. Q_v is the set of verification queries issued by the same adversary. Then, the probability that this adversary can successfully create a forgery is bounded in the following way:*

$$\text{Adv}(\mathcal{A}_{MAC}^{ws\text{-}MAC(), Q_s, Q_v}) \leq \frac{m \cdot |Q_v|}{q} + (2)$$
$$|Q_v| \cdot (\text{Adv}_{|Q|_{00}}^{E_{00}()} + \text{Adv}_{|Q|_{01+1}}^{E_{01}()} + \text{Adv}_{|Q|_{10+n}}^{E_{10}()})$$

where $|Q|_{00} = \frac{n \cdot m \cdot w_e}{w_c} \cdot |Q_s|$, $|Q|_{01} = |Q_s| + |Q_v|$, $|Q|_{10} = n \cdot (|Q_s| + |Q_v|)$, and q is the prime found in the definition of Algorithms 2 and 3.

The security level depends on the prime number q , number of queries served $|Q_s|$, $|Q_v|$, the dimensions of the matrix in terms of number of rows n and columns m , and the security of block ciphers $E_{00}()$, $E_{01}()$ and $E_{10}()$. We choose a prime number q that is the largest prime number in 2^{w_t} . For example, we use $w_t = 127$ and $q = 2^{127} - 1$, considering both security and performance. If we consider a 1024-dimension matrix row, we can serve 2^{53} queries without changing key, while maintaining a security level higher than 64 bits.

V. ARCHITECTURAL SUPPORT FOR SECNDP

The SecNDP scheme can be applied to any TEE (CPU, GPU, ASIC accelerator, etc.) and work with any untrusted near-memory or near-storage processing hardware. In this section, we describe a SecNDP architecture and design for a computing platform supported by a TEE-enabled processor with near-memory processing.

Baseline NDP Architecture: Figure 5 depicts the components supporting NDP in a dual in-line memory module (DIMM). NDP architecture support (blue boxes) comprises of an NDP protocol (i.e., NDP commands), CPU ISA extensions for issuing NDP packets, and NDP PUs on the memory side. NDP PUs can compute and access their memory region in parallel. For example, the Rank-NDP PUs run in parallel to access the memory within their DRAM rank simultaneously and perform computations to obtain the intermediate results. Therefore, the peak compute and memory throughput of the NDP grows with the number of ranks. Each NDP PU contains registers to hold its intermediate results. Multiple register allow multiple NDP operations to overlap without sending intermediate results back to a CPU. For workloads that need to store a number of intermediate results simultaneously, the number of NDP PU registers can become the bottleneck and more registers can improve performance.

ISA extensions for NDP: The processor issues special instructions that offload NDP packets to the memory controller, which then dispatches NDP commands (Figure 5) to the NDP PUs. There are two types of NDP instructions (and the corresponding commands): $NDPInst$ that controls the arithmetic computation, and $NDPLd$ that loads the value in an NDP register from and to the processor. $NDPInst$ has all the operands for issuing an NDP command, including a data address, the operation Op to be carried out, vector size $vsize$, data size $dsize$, an immediate operand value Imm , and source/destination register IDs $RegID$. Figure 5 also shows the value in an NDP command for vector-matrix multiplication $a \times P$. With this instruction, an NDP command will be issued to let an NDP PU first multiply each element

within the same DRAM row, saving memory access time and energy. (2) *Ver-sep*: allocating tags separately from the data. At the boot-up time or TEE loading time, a system can allocate a designated physical address region to store tags. The advantage of this option is that the software binary layout does not need to be changed for verification. (3) *Ver-ECC*: storing the tags in an ECC chip, and storing the ECC bits in separate memory space [63]. Both data and its tag are fetched in one memory access, while ECC bits are fetched only when corruption is detected. The advantage is that it co-locates the tag with data and fully utilizes the bandwidth of commercial hardware. The disadvantage is that the ECC chip has a fixed capacity and is not flexible for different vector and tag sizes. Also, the SecNDP scheme only verifies the final computation result over a potentially large memory region. If the verification fails, all the ECC of related memory regions need to be checked.

When an NDP operation is to be verified, the NDP needs to compute on the tags. The operation on the data vector ($a_i \times C_i$) will also apply to the tag ($a_i \times C_{T_i}$), as shown in Figure 4(a). There are two possible designs. The SecNDP Engine can issue an extra NDP instruction to compute the tag. Another possible design is to extend each operation. Instead of operation on a vector ($a \times C_i$), an operation on a vector and a tag $a \times [C_i | C_{T_i}]$ is conducted. The second design needs to change the NDP PU logic to have extended registers. Note that the computation for tags is in a prime field. If we choose $q = 2^{127} - 1$, the computation is similar to standard integer arithmetic, but has extra logic when an overflow is incurred [13].

E. Implementation of SecNDP Instructions

1) *ArithEnc*: The AES engines use *paddr* and *v* to generate the OTPs. Each plaintext data value is then subtracted by its OTP bits (following Alg. 1) and written back to memory like a cache line flush. If the verification bit is set, a MAC (Alg. 2) is computed in the verification engine, and the encryption engine uses *paddr*, *m*, *v* to encrypt the tag (Alg. 3).

2) *SecNDPInst*: The memory controller passes the operands to the SecNDP Engine and also issues NDP commands. As shown in Figure 4(a), the OTP PU and the NDP PU perform the same computation on the OTP and the ciphertext, respectively. The intermediate results will be in the corresponding registers in the OTP PU and the NDP PU. If verification is needed, extra computation on the tag is performed as described in Section V-D.

3) *SecNDPLd*: On the *SecNDPLd* instruction, the value in the desired NDP PU register is loaded to the *resp. buf.* in the processor (see Figure 5). The value of the corresponding register in the OTP PU register is loaded to *dec. buf.* and added to the value from NDP. Only one adder is on the critical path when the encrypted result is returned to the processor core. If verification is needed,

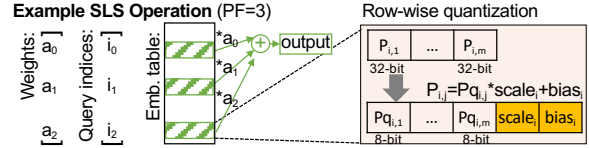


Figure 6. (Left) Embedding table lookup (SLS operation) in deep learning recommendation model. (Right) Row-wise quantization to reduce memory footprint.

the data result is further passed to the Verification Engine to compute the checksum. Meanwhile, the tag is obtained by adding the intermediate tag in the OTP PU and the NDP PU. If the tag matches the checksum of the data result, the verification passes. Otherwise, the verification fails and an interrupt will be triggered. The verification can be in the critical path (1–2 cycle), or can be speculated [48].

VI. EVALUATION METHODOLOGY

A. Workloads

We use two categories of representative NDP use cases. The first one is deep learning recommendation inference [28], [58], which exhibits sparse and irregular memory access patterns. The second one is data analytics for biomedical applications, which perform accesses over contiguous memory regions. A pooling factor (*PF*) denotes the level of data aggregation, *i.e.*, the ratio of the raw data size to the size of the computed result.

(1) Deep Learning Recommendation Inference: Recommendation models, such as Deep Learning Recommendation Model (DLRM) [58] are structured to take advantage of both continuous and categorical features of individual users. In addition to continuous feature processing using fully-connected (FC) layers commonly found in deep learning, the categorical features are captured by large embedding tables with sparse lookup and pooling operations.

Embedding tables are organized as a set of potentially millions of vectors of dimension m . An embedding table lookup operation (*i.e.*, SparseLengthsSum or SparseLengthsWeightedSum (SLS) operation in Caffe2 [3]) performs a weighted summation for a set of vectors. As illustrated in the left part of Figure 6, an SLS query consists of a list of indices $[i_0, i_1, \dots, i_{PF-1}]$ and a weight vector a of dimension PF . The result of the pooling is a vector of dimension m and each element of the vector is $res_j = \sum_{k=0}^{PF-1} a_{i_k} \times P_{i_k,j}$. The embedding operations in recommendation models often access irregular indices from large tables, impeding the memory performance that further limits the overall service throughput. Previous studies showed that NDP can significantly speed up embedding lookups in recommendation models [36], [76].

In our evaluation, we offload the embedding table lookup (*i.e.*, SLS operations) to NDP and the rest of the model runs on a CPU. We consider the recommendation models as the

Table I
PARAMETERS OF DLRM MODELS IN EVALUATION

	bottom FC	top FC	# Emb.	total Emb. size
RMC1-small	256-128-32	256-64-1	8	1GB
RMC1-large	256-128-32	256-64-1	12	1.5GB
RMC2-small	256-128-32	256-128-1	24	3GB
RMC2-large	256-128-32	256-128-1	64	8GB

service provider’s intellectual property (IP), while the values containing users’ private information in the embedding tables as the data that need to be protected.

Memory footprint optimization by quantization. For efficiency, the values in the embedding tables are often quantized to a smaller bit width [21]. For example, in row-wise quantization, a vector of 32-bit values can be quantized into a vector of 8-bit integers with a scale and bias per row (Figure 6 right). When the vector is queried, each value in the vector is multiplied by the scale and then the bias is added to recover the original value ($P_{i,j} = Pq_{i,j} \times scale_i + bias_i$). Proper quantization schemes reduce memory footprint while keeping the model accuracy [21].

However, when the row-wise quantization is applied, for each computation there is an extra multiplication with scale ($Pq_{i,j} \times scale_i$), making computation over ciphertext less efficient. Thus, we propose and evaluate table-wise and column-wise quantization schemes, where the scale and the bias are assigned on a per-table or a per-column basis. With table-wise and column-wise quantization, the SLS operation can be first performed without per-row scale, *i.e.*, $resq_j = \sum_{k=0}^{PF-1} a_{ik} \times Pq_{ik,j}$. In the last step, the per-table or per-column scale and bias are used to get the final result, *i.e.*, $res_j = resq_j \times scale_j + bias_j$. The total size of the per-column or the per-table scale and the bias is much smaller than the table size. With the table-wise or column-wise quantization, the vector-matrix multiplication can be directly applied to quantized values. For performance evaluation, we assume the scale and bias for quantization can be cached in the processor, and thus, the table-wise and column-wise quantization have the same performance.

Recommendation Model Parameters. For performance evaluation, we use DLRM models with the representative model parameters in Table I. Each embedding table row has $m = 32$ elements. We also consider a different quantization scheme that quantizes 32-bit values into 8-bit values (*i.e.*, 2 cache lines into about 0.5 cache line per vector.) We use a randomly-generated query trace with $PF = 40$ and 80, and a query trace from a production model with a pooling factor PF ranging from 50 to 100. For encryption and replay attack prevention, each embedding table uses a version number, and the enclave software manages at most 64 version numbers. For model accuracy evaluation, we use a production-scale recommendation model consisting of hundreds of embedding tables with production data set.

(2) Medical Data Analytics. The second use case we consider here is data analytics over a private data set. To study whether a certain disease is related to certain genes,

Table II
SIMULATION PARAMETERS AND CONFIGURATIONS

DRAM Parameters
DDR4-2400MHz, rank_size=8GB, tRC=55, tRCD=16, tCL=16, tRP=16, tBL=4, tCCD S=4, tCCD L=6, tRRD S=4, tRRD L=6, tFAW=26
AES Encryption Engine Parameter [22]
Throughput = 111.3Gbps, <i>i.e.</i> , 1.15ns per 128-bit block.

statistical hypothesis tests need to be performed. Consider a data set containing the expression level for $m = 10000$ genes of $n = 500,000$ patients (e.g., [16]). This is used to compute the test statistics (e.g., p-value of t-test [71]), and the summation (or average) of the gene expression level of patients and non-patients.

The data set contains sensitive medical information (e.g., the gene expression level) and needs to be protected. We store the data set in the memory after encryption. When researchers want to study a certain disease, they query the data set by giving a list of patient IDs in the data set, and let the NDP unit compute the summation. A query is comprised of a list of patient IDs which are used for aggregating their gene expression level. Usually the queried patient IDs are not sparse.

Database Parameters: In performance simulation, we consider a database with $m = 1024$ genes and conduct summation over $PF = 10,000$ patients (40MB in total).

B. Evaluation Setup

SecNDP Performance Simulation. Based on the simulation framework in [36], we built a cycle-level simulation framework with the following components: (1) a physical addresses mapping module, (2) an NDP packet generator, (3) an encryption engine, and (4) an NDP DIMM consisting of DRAM devices, arithmetic units, and control logic. We use Ramulator [40], a cycle-level DRAM simulator, for the DRAM devices. Table II summarizes the parameters and configurations.

During the simulation, we emulate the packet generation and scheduling steps taken by the software stack and the memory controller. First, we apply a standard page mapping method [26] to generate the physical addresses from a trace of embedding lookups by assuming that the OS randomly selects free physical pages for each logical page frame. This physical address trace is then fed to Ramulator.

We implemented a cycle-level NDP module on top of Ramulator including the logic in the NDP PU. The NDP_rank parameter denotes the number of Rank-NDP PUs in the system, and the parameter NDP_reg denotes the number of registers per NDP PU. To evaluate the NDP latency, the packet generator divides the physical memory requests into packets of NDP commands that are then sent to the cycle-level simulator. NDP activates all ranks under the memory channel. For every NDP packet, the NDP commands are dispatched to the parallel ranks and the latency is bounded by the slowest rank. The total latency of NDP also includes

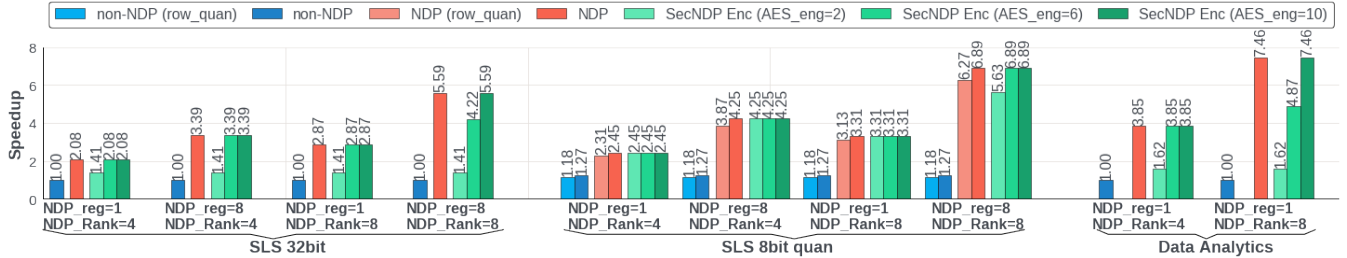


Figure 7. Performance of unprotected non-NDP baseline (in blue), unprotected NDP (in red), and SecNDP-Enc with different numbers of AES engines (in green). For SLS 8-bit quantization, (row_quan) denotes the result of the row-wise quantization scheme, and the other bars are for both the column-wise and table-wise quantization schemes⁵. With a sufficient number of AES engines, encryption in SecNDP shows the same speedup as unprotected NDP.

the DRAM cycles during initialization to configure memory-mapped control registers and a cycle in the final stage to transfer the sum/partial-sum to the processor using NDPLd.

To evaluate the throughput of the AES engine and the OTP PU, we use the performance number of a fully pipelined AES design [22]. We assume the addition and multiplication on the counter block are pipelined cycle-by-cycle after AES encryption. Combining the off-chip NDP latency and the encryption engine throughput, we estimate the throughput of SecNDP. The final throughput is the smaller one of the NDP throughput and the OTP throughput. **Whole System (TEE + SecNDP) Evaluation.** In some cases, only part of the computation will be offloaded to SecNDP (NDP portion). The remaining part will be executed in the CPU TEE (CPU portion). To estimate the performance slowdown of the CPU portion and to establish a CPU TEE performance reference, we measure the execution time on two Intel machines with SGX enclaves [4]. One is an Intel Xeon E-2288G CoffeeLake (CFL) CPU, with 168MB SGX EPC, 16MB L3, and 32GB DRAM. The other one is an Intel Xeon Platinum 8370C IceLake (ICL) CPU, with 96GB SGX EPC, 48MB L3, and 192GB DRAM. When the workload fits in caches (e.g., the CPU portion of DLRM), SGX ICL has about 5% slowdown. To evaluate the end-to-end performance, we ran the whole model on the ICL machine to obtain the execution time breakdown between the CPU portion and the NDP portion, and use the speedup (or slowdown) of each portion to evaluate the execution time. **Arithmetic Precision and Application Accuracy.** We apply a quantization scheme to a production-scale model and use a production dataset with 40K samples to evaluate the model accuracy. The accuracy is presented using *Logloss* [21], a widely used loss function for prediction models.

Power and Area. We use the memory trace from our performance simulation setup and DRAMPower tool [18] to estimate the DRAM chip energy. We use CACTI-IO [34]

⁵As discussed in Section VI-A-(1), because we cache the scale and bias for quantization in the processor, the column-wise and the table-wise quantization schemes have the same speedups. Hence, we display one bar for both the column-wise and table-wise quantization in the middle group of SLS 8bit quan. The (row_quan) bars represent the original row-wise quantization scheme in the baseline and unprotected NDP scenarios.

Table III
SECNDP SPEEDUP AGAINST UNSECURED BASELINE AND SGX.

	RMC1-small	RMC1-large	RMC2-small	RMC2-large	Data Analytics
unprotected non-NDP	1x	1x	1x	1x	1x
unprotected NDP	2.46x	3.11x	4.05x	4.44x	7.46x
SGX-CFL	0.0038x	0.0037x	N/A	N/A	0.1738x
SGX-ICL (no int. tree)	0.59x	0.60x	N/A	N/A	0.57x
SecNDP	2.36x	3.02x	3.95x	4.33x	7.46x

to evaluate the energy between the DRAM chip and the NDP PU (located inside DIMM’s buffer chip) and the energy of DIMM IO. To estimate the energy and the area of the SecNDP engine, we refer to an AES design [22] and use model in [66] for the OTP PU and the verification engine at 45nm process node.

VII. EVALUATION RESULTS

In this section, we evaluate the performance (in execution time), precision, energy, and area of the SecNDP architecture. The results show the performance and energy consumption of the SecNDP approach and that of unprotected NDP. The precision loss in SecNDP shows negligible impact on the accuracy of the recommendation model. The area overhead of the SecNDP engine is also small.

A. Performance Evaluation

Overall End-to-End Performance. Table III illustrates the performance of SecNDP compared with the non-NDP baseline and Intel SGX. By leveraging eight NDP PUs, SecNDP demonstrates 2.3x to 4.3x speedup for end-to-end DLRM models with *batch size*=256 and 7.46x for the medical data analytics. The performance of SecNDP is close to that of an unprotected NDP as also shown in the table. Meanwhile, using SGX shows considerable slowdown (Table III). This is because the working set sizes of the workloads do not fit in on-chip caches. Due to the malloc size limit by the current SGX library, we could only run RMC1 in SGX. We observe 6x–300x slowdown for the CFL SGX enclave⁶, and 1.8x – 2.6x slowdown for the ICL SGX enclave⁷. For SecNDP,

⁶CFL processors rely on an integrity tree to protect against replay attacks, and thus only support a limited number of protected memory pages, causing frequent page swapping in this case.

⁷ICL processors do not have integrity tree for replay attack prevention.

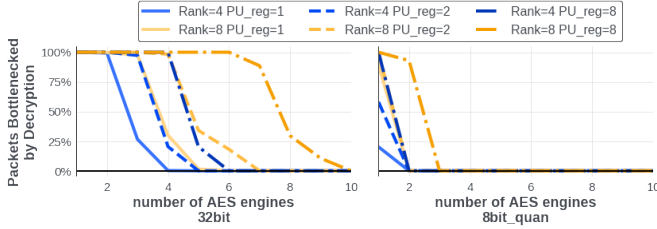


Figure 8. Percentage of NDP packets for SLS operations bottlenecked by decryption bandwidth for confidentiality protection of SecNDP.

Table III shows the performance with the verification scheme (*Ver-ECC*) that stores tags in the ECC chip. We will present the performance of unprotected NDP and SecNDP with various verification schemes (introduced in Section V-D) in more detail next. Overall, SecNDP demonstrates significant performance improvement over a CPU TEE without NDP.

NDP Performance - Unprotected. The red bars in Figure 7 summarize the performance of the SLS operations (NDP portion of DLRM) using different quantization schemes and that of the data analytics workload across different NDP settings (NDP_rank , NDP_reg). For SLS operations, quantization provides 17–27% speedup for all row-wise (row_quan), column-wise, and table-wise quantization in both the NDP and non-NDP settings. This is because the embedding tables are smaller after quantization. Further, with more NDP_rank and NDP_reg , NDP tends to have higher speedup. More NDP_rank allows more parallel memory accesses. More NDP_reg allows more intermediate SLS results, making the workload among the NDP PU more evenly distributed. For $NDP_rank=8$ and $NDP_reg=8$, the speedup reaches 5.59x without quantization, and 6.89x with quantization. The data analytics workload exhibits more regular memory access patterns that are distributed evenly across all the NDP_rank PUs. Thus, it results in higher performance speedup (up to 7.46x) than irregular SLS operations. Also, since there is only one resulting sum, more NDP_reg does not help further for the data analytics workload.

NDP Performance - Encryption-Only. We evaluate the performance of SecNDP over different numbers of AES engines (Figure 7). When there is only a small number of AES engines, the decryption becomes the performance bottleneck. As the number of AES engines increases, the speedup reaches that of the unprotected NDP in all settings, indicating that the performance bottleneck eventually shifts to the memory bandwidth.

Figure 8 shows the percentage of NDP packets that is bottlenecked by the AES bandwidth for SLS operations. If the workload has enough parallelism and no NDP PU is idle, more NDP_rank require more AES engines to match the NDP throughput. For example, when $NDP_rank=8$, we need ten AES engines to match the memory throughput in the burst mode for the system setting in Table II. However, in practice, applications may not be able to fully utilize the

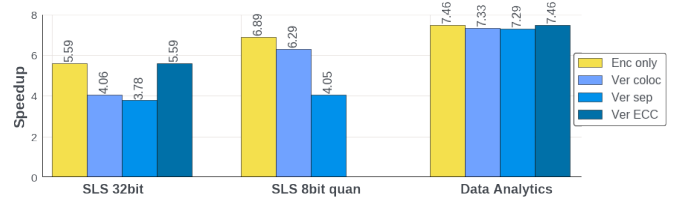


Figure 9. Speedup of various SecNDP encryption and verification schemes with $NDP_rank=8$, $NDP_reg=8$.

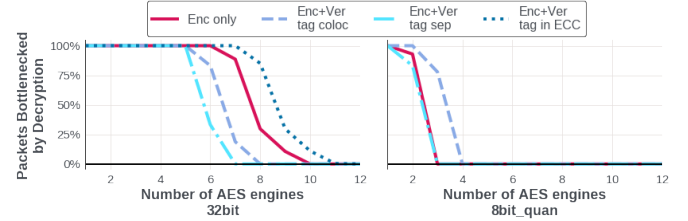


Figure 10. Percentage of NDP packets for SLS operations bottlenecked by decryption bandwidth with $NDP_rank=8$, $NDP_reg=8$.

peak memory bandwidth and less AES engines may be sufficient. For example, Figure 8 shows that eight AES engines are enough to match the NDP throughput for roughly 70% of the NDP packets for SLS operations without quantization. Quantization can also significantly reduce the number of AES engines needed; with quantization, only about one third of the AES engines are needed. This is because less OTP is required for the decryption.

NDP Performance - Encryption+Verification. We evaluate three design options for SecNDP’s verification tags (*Ver-coloco*, *Ver-sep*, and *Ver-ECC* in Section V-D), and assume the NDP PUs have enough bandwidth to process tags. We use a 128-bit tag for each vector. Figure 9 shows the performance results of $NDP_rank=8$ and $NDP_reg=8$ with twelve AES engines. For the data analytics workload, because the row vector size m is greater than that of SLS, and the 128-bit tag is relatively small compared to data, the verification overhead is small. Figure 10 presents the percentage of the NDP packets bottlenecked by the decryption operation for SLS operations.

As shown in Figure 9, without quantization, *Ver-coloco* and *Ver-sep* show lower speedup because of the extra memory accesses for the tags. *Ver-sep* shows more performance degradation because tags and data are not co-located, in other words, an additional DRAM row buffer activation is required. *Ver-ECC* has the same speedup as the encryption-only (*Enc-only*) case as no extra DRAM access is required. As shown in the left part of Figure 10, *Ver-ECC* needs more AES engines to match the bandwidth requirement.

With quantization, the embedding vector size becomes shorter than a cache line and the corresponding tags cannot fit in the ECC chip; *Ver-ECC* does not work. Hence, we only show the results for *Ver-coloco* and *Ver-sep* in Figure 9. With *Ver-coloco*, one cache line fetch can retrieve both data and tag, and thus the performance is close to *Enc-only*. But *Ver-coloco*

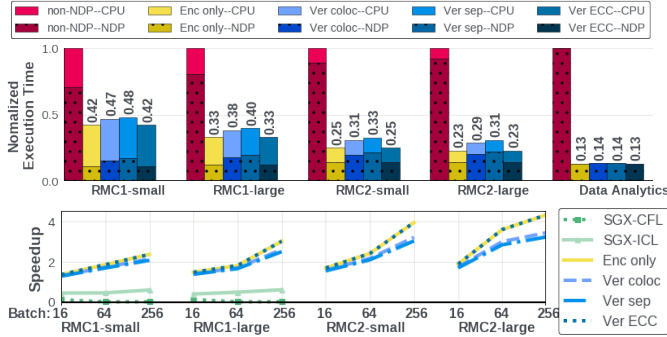


Figure 11. (Top) Normalized execution time in SecNDP with $NDP_rank=8$. The breakdown execution time of NDP portion (–NDP) and the CPU TEE portion (–CPU) of the workloads is shown in the stacked bar. (Bottom) Recommendation inference speedup for different batch sizes in SecNDP.

Table IV
ACCURACY OF DIFFERENT QUANTIZATION SCHEMES

	LogLoss	LogLoss degradation
32-bit floating point	0.64013	0
32-bit fixed point	0.64013	$-3.6 * 10^{-10}$
table-wise quantization (8-bit)	0.64059	0.07%
column-wise quantization (8-bit)	0.64027	0.02%

still cannot reach the performance of *Enc-only*, because with the tag present, the data is not aligned with the cache line boundary. In some cases, two contiguous cache lines are still needed. Compared to *Enc-only*, *Ver-coloc* requires more AES engines to decrypt the tag in parallel with the data. *Ver-sep* leads to about 40% performance degradation over *Enc-only*, because two separate cache lines are accessed instead of one. Even though *Ver-sep* does not change an application’s memory layout and can support a variety of tag sizes, it has the worst performance. Co-locating tags and data either in ECC or in the memory has better performance but requires a more complicated system implementation.

End-to-End Execution Time Breakdown. Figure 11 presents the breakdown of the end-to-end execution time and how the performance speedup scales with different batch sizes. Here, we present the speedup with $NDP_rank=8$, $NDP_reg=8$, and $PF=80$ for SLS operations without quantization. For $batch\ size=256$, compared with the insecure non-NDP baseline, SecNDP achieves an end-to-end model inference speedup, ranging from 2.3x to 4.3x. The results include the slowdown of CPU portion in TEE. SecNDP provides higher speedup for larger batch sizes, while SGX does not scale with batch sizes.

B. Impact of Arithmetic Precision on Application Accuracy

Like in MPC, we only support integer and fixed point operations in SecNDP. To assess the impact of data type precision, we evaluate the accuracy impact from the precision change [32]. We perform the model accuracy evaluation using production-level recommendation models [36]. As shown in Table IV, using 32-bit fixed-point values for the embedding has a negligible impact on the Logloss. Further-

Table V
MEMORY ENERGY CONSUMPTION OF SECNDP (PJ/BIT)

	DIMM	DIMM IO	SecNDP Engine	Normd. Mem. Energy (PF=80)
unprotected non-NDP	$27.42 \times PF$	$7.3 \times PF$	0	100%
unprotected NDP	$27.42 \times PF$	7.3	0	79.2%
non-NDP Enc	$27.42 \times PF$	$7.3 \times PF$	$0.5 \times PF$	101.5%
SecNDP Enc	$27.42 \times PF$	7.3	$0.9 \times PF$	81.83%
SecNDP Enc+ver	$30.85 \times PF$	8.2	$1.01 \times PF + 1.72$	92.09%

more, using 8-bit table-wise or column-wise quantization leads to similar, negligible Logloss degradation ($<0.07\%$). The accuracy results suggest that SecNDP can provide significant speedup with security guarantee, while maintaining the desired model accuracy requirement.

C. Energy and Area Overhead

In addition to performance improvement, NDP also saves energy by reducing the amount of data transfer between the processor and DIMM. Table V summarizes the energy consumption of memory (including IO) and the SecNDP engine. The *non-NDP Enc* row shows the energy consumption of a TEE without NDP. Although SecNDP requires extra computation in OTP PUs, the energy consumption of OTP PUs is insignificant compared to the overall memory power savings. When $PF=80$, SecNDP saves memory system energy by 18% with encryption only and by 8% with verification. In addition, there will be additional energy savings on the processor cache hierarchy, because less data is moved in and out of the caches. The area overhead of SecNDP is estimated to be 1.625 mm^2 at 45nm node if we use 10 AES engines which match the throughput of OTP PUs and the verification engine. The energy/area overheads can be further reduced with more advanced process nodes.

VIII. RELATED WORK

Many cryptographic schemes have been proposed to protect the confidentiality of sensitive workloads on untrusted platforms, such as HE [15], [19], [23], [53], [59], [60] and MPC [20], [55]. However, HE incurs about several orders of magnitude performance overheads [60]. Thus, deploying HE in untrusted NDP does not outperform the TEE baseline without NDP. MPC has better performance than HE, however, the security of MPC relies on the assumption that the untrusted nodes do not collude, which is not realistic in many system settings. To further verify the correctness of computations results, non-interactive verifiable computation [25], [50] and probabilistically checkable proofs [12], [65] are proposed, but their computation complexity is still too high to be practical for secure outsourcing.

Hardware-based TEE is an alternative solution to protect both confidentiality and integrity of workloads. Analyzing the memory protection of Intel SGX, Vessels [39] showed that deep learning workloads suffer from significant performance degradation, and proposed data movement optimizations. To improve the performance of TEE for data-intensive workloads, untrusted accelerators can be leveraged [31],

[74]. Slalom [74] proposed using arithmetic secret sharing to offload secure computations from a TEE to untrusted GPUs. However, the TEE still needs to store its share of secret in memory and pre-compute the results in an offline phase. Thus, Slalom moves computation from online to offline, but does not reduce computation or memory usage. DarKnight [31] proposed a blinding scheme to offload batched deep learning workloads to untrusted GPU. However, DarKnight cannot protect weights in the accelerator. Another solution is to also trust the accelerators and include the accelerators in the TEE. For memory, previous studies [5], [9] proposed protocols to encrypt and obfuscate the traffic on the memory bus. However, InvisiMem [5] still incurs significant performance, energy and memory space overhead. In addition, such solution not only requires trusting multiple hardware vendors, but also requires coordination and standardization among the vendors. SecNDP is the first work to demonstrate how untrusted off-the-shelf NDP units can be used for secure computation.

IX. CONCLUSION

The recent progress in near-data processing generates a lot of interest in using NDP to alleviate the memory bandwidth bottleneck for data-intensive applications. However, there is a lack of feasible techniques that protect the confidentiality of off-chip data while taking advantage of NDP. In this paper, we present SecNDP, a lightweight encryption and verification scheme that supports NDP over ciphertext and verifies the correctness of NDP results. With a sufficient number of AES engines, our evaluation on two workloads shows that SecNDP can match the speedup delivered by unprotected NDP. Our energy evaluation shows SecNDP saves memory energy by reducing the data transfer on the memory bus. The SecNDP scheme enables a TEE in the presence of untrusted memory to leverage the performance and energy benefits of NDP securely.

ACKNOWLEDGMENTS

The authors would like to thank Muhammad Umar, Henry Wang, Shankaran Gnanashanmugam, Jihang Li, Yuchen Hao, and Haixin Liu for their help in evaluating recommendation system in Intel SGX, and thank Brian Knott, Hao Chen, Chuan Guo for their suggestions. The authors would also like to thank the anonymous reviewers for their insightful comments and suggestions. Liu Ke and Xuan Zhang were partially supported by NSF CCF-1942900.

REFERENCES

- [1] “AMD Secure Encrypted Virtualization (SEV),” <https://developer.amd.com/sev/>.
- [2] “Arm TrustZone Technology,” <https://developer.arm.com/ip-products/security-ip/trustzone>.
- [3] “Caffe2,” <https://caffe2.ai/>.

- [4] “Intel Software Guard Extensions(SGX),” <https://www.intel.com/content/www/us/en/architecture-and-technology/software-guard-extensions.html>.
- [5] S. Aga and S. Narayanasamy, “InvisiMem: Smart memory defenses for memory bus side channel,” in *International Symposium on Computer Architecture (ISCA)*, 2017, p. 94–106.
- [6] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi, “A scalable processing-in-memory accelerator for parallel graph processing,” in *International Symposium on Computer Architecture (ISCA)*, 2015, p. 105–117.
- [7] J. Ahn, S. Yoo, O. Mutlu, and K. Choi, “PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture,” in *International Symposium on Computer Architecture (ISCA)*, 2015, pp. 336–348.
- [8] H. Asghari-Moghaddam, Y. H. Son, J. H. Ahn, and N. S. Kim, “Chameleon: Versatile and Practical Near-DRAM Acceleration Architecture for Large Memory Systems,” in *International Symposium on Microarchitecture (MICRO)*, 2016, pp. 1–13.
- [9] A. Awad, Y. Wang, D. Shands, and Y. Solihin, “ObfusMem: A low-overhead access obfuscation for trusted memories,” in *International Symposium on Computer Architecture (ISCA)*, 2017, p. 107–119.
- [10] Bahar Asgari, Ramyad Hadidi, Jiashen Cao, Da Eun Shim, Sung-Kyu Lim, Hyesoon Kim, “FAFNIR: Accelerating Sparse Gathering by Using Efficient Near-Memory Intelligent Reduction,” in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2021, pp. 908–920.
- [11] S. Banerjee, P. Ramrakhiani, S. Wei, and M. Tiwari, “SESAME: Software defined enclaves to secure inference accelerators with multi-tenant execution,” *arXiv preprint arXiv:2007.06751*, 2020.
- [12] M. Bellare, S. Goldwasser, C. Lund, and A. Russell, “Efficient probabilistically checkable proofs and applications to approximations,” in *Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*, 1993, pp. 294–304.
- [13] D. J. Bernstein, “Floating-point arithmetic and message authentication,” 2000. [Online]. Available: <http://cr.yp.to/papers.html#hash127>.
- [14] T. Bourgeat, I. Lebedev, A. Wright, S. Zhang, Arvind, and S. Devadas, “MI6: Secure Enclaves in a Speculative Out-of-Order Processor,” in *International Symposium on Microarchitecture (MICRO)*, 2019, p. 42–56.
- [15] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, “Fully homomorphic encryption without bootstrapping,” in *Innovations in Theoretical Computer Science (ITCS)*, 2012, pp. 309–325.
- [16] C. Bycroft, C. Freeman, D. Petkova, G. Band, L. T. Elliott, K. Sharp, A. Motyer, D. Vukcevic, O. Delaneau, J. O’Connell *et al.*, “The UK biobank resource with deep phenotyping and genomic data,” *Nature*, vol. 562, no. 7726, pp. 203–209, 2018.
- [17] R. Cammarota, I. Banerjee, and O. Rosenberg, “Machine Learning IP Protection,” in *Proceedings of the International Conference on Computer-Aided Design (ICCAD)*, 2018.
- [18] K. Chandrasekar, B. Akesson, and K. Goossens, “Improved power modeling of ddr sdrams,” in *2011 14th Euromicro Conference on Digital System Design*, 2011, pp. 99–108.
- [19] J. H. Cheon, A. Kim, M. Kim, and Y. Song, “Homomorphic encryption for arithmetic of approximate numbers,” in *Annual International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT)*, 2017.

- [20] I. Damgård, V. Pastro, N. Smart, and S. Zakarias, "Multiparty computation from somewhat homomorphic encryption," in *Annual Cryptology Conference (CRYPTO)*, 2012, pp. 643–662.
- [21] Z. S. Deng, J. Park, P. T. P. Tang, H. Liu, J. Yang, H. Yuen, J. Huang, D. S. Khudia, X. Wei, E. Wen, D. Choudhary, R. Krishnamoorthi, C.-J. Wu, N. Satish, C. Kim, M. Naumov, S. Naghshineh, and M. Smelyanskiy, "Low-precision hardware architectures meet recommendation model inference at scale," *IEEE Micro*, 2021.
- [22] P.-K. Dong, H. K. Nguyen, and X.-T. Tran, "A 45nm high-throughput and low latency aes encryption for real-time applications," in *2019 19th International Symposium on Communications and Information Technologies (ISCIT)*, 2019, pp. 196–200.
- [23] J. Fan and F. Vercauteren, "Somewhat practical fully homomorphic encryption." *IACR Cryptol. ePrint Arch.*, vol. 2012, p. 144, 2012.
- [24] A. Farmahini-Farahani, J. H. Ahn, K. Morrow, and N. S. Kim, "NDA: Near-DRAM Acceleration Architecture Leveraging Commodity DRAM Devices and Standard Memory Modules," in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2015, pp. 283–295.
- [25] R. Gennaro, C. Gentry, and B. Parno, "Non-interactive verifiable computing: Outsourcing computation to untrusted workers," in *Annual Cryptology Conference (CRYPTO)*. Springer, 2010, pp. 465–482.
- [26] M. Gorman, "Understanding the Linux virtual memory manager," 2004.
- [27] C. Guo, A. Hannun, B. Knott, L. van der Maaten, M. Tygert, and R. Zhu, "Secure multiparty computations in floating-point arithmetic," *arXiv preprint arXiv:2001.03192*, 2020.
- [28] U. Gupta, C.-J. Wu, X. Wang, M. Naumov, B. Reagen, D. Brooks, B. Cottel, K. Hazelwood, M. Hempstead, B. Jia, H.-H. S. Lee, A. Malevich, D. Mudigere, M. Smelyanskiy, L. Xiong, and X. Zhang, "The architectural implications of facebook's dnn-based personalized recommendation," in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2020, pp. 488–501.
- [29] J. A. Halderman, S. D. Schoen, N. Heninger, W. Clarkson, W. Paul, J. A. Calandrino, A. J. Feldman, J. Appelbaum, and E. W. Felten, "Lest We Remember: Cold-Boot Attacks on Encryption Keys," *Commun. ACM*, vol. 52, no. 5, p. 91–98, May 2009.
- [30] S. Halevi and H. Krawczyk, "MMH: Software message authentication in the gbit/second rates," in *Fast Software Encryption*, 1997, pp. 172–189.
- [31] H. Hashemi, Y. Wang, and M. Annavaram, "DarKnight: An accelerated framework for privacy and integrity preserving deep learning using trusted hardware," in *International Symposium on Microarchitecture (MICRO)*, 2021, pp. 212–224.
- [32] N. J. Higham, *Accuracy and Stability of Numerical Algorithms*, 2nd ed. Society for Industrial and Applied Mathematics, 2002.
- [33] W. Hua, M. Umar, Z. Zhang, and G. E. Suh, "GuardNN: Secure DNN Accelerator for Privacy-Preserving Deep Learning," *arXiv preprint arXiv:2008.11632*, 2020.
- [34] N. P. Jouppi, A. B. Kahng, N. Muralimanohar, and V. Srinivas, "Cacti-io: Cacti with off-chip power-area-timing models," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, pp. 1254–1267, 2014.
- [35] Y. Kang, Y.-s. Kee, E. L. Miller, and C. Park, "Enabling cost-effective data processing with smart SSD," in *Symposium on Mass Storage Systems and Technologies (MSST)*, 2013.
- [36] L. Ke, U. Gupta, B. Y. Cho, D. Brooks, V. Chandra, U. Diril, A. Firoozshahian, K. Hazelwood, B. Jia, H.-H. S. Lee, M. Li, B. Maher, D. Mudigere, M. Naumov, M. Schatz, M. Smelyanskiy, X. Wang, B. Reagen, C.-J. Wu, M. Hempstead, and X. Zhang, "RecNMP: Accelerating personalized recommendation with near-memory processing," in *International Symposium on Computer Architecture (ISCA)*, 2020, pp. 790–803.
- [37] L. Ke, X. Zhang, J. So, J.-G. Lee, S.-H. Kang, S. Lee, S. Han, Y. Cho, J. H. Kim, Y. Kwon, K. Kim, J. Jung, I. Yun, S. J. Park, H. Park, J. Song, J. Cho, K. Sohn, N. S. Kim, and H.-H. S. Lee, "Near-Memory Processing in Action: Accelerating Personalized Recommendation with AxDIMM," *IEEE Micro*, 2022.
- [38] D. Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay, "Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory," in *International Symposium on Computer Architecture (ISCA)*, 2016, pp. 380–392.
- [39] K. Kim, C. H. Kim, J. J. Rhee, X. Yu, H. Chen, D. J. Tian, and B. Lee, "Vessels: Efficient and scalable deep learning prediction on trusted processors," in *Symposium on Cloud Computing (SoCC)*, 2020, pp. 462–476.
- [40] Y. Kim, W. Yang, and O. Mutlu, "Ramulator: A fast and extensible DRAM simulator," in *IEEE Computer architecture letters*, vol. 15, no. 1, 2015, pp. 45–49.
- [41] B. Knott, S. Venkataraman, A. Hannun, S. Sengupta, M. Ibrahim, and L. van der Maaten, "CrypTen: Secure multiparty computation meets machine learning," in *Proceedings of the NeurIPS Workshop on Privacy-Preserving Machine Learning*, 2020.
- [42] T. Kohno, J. Viega, and D. Whiting, "CWC: A high-performance conventional authenticated encryption mode," in *International Workshop on Fast Software Encryption*. Springer, 2004, pp. 408–426.
- [43] Y.-C. Kwon, S. H. Lee, J. Lee, S.-H. Kwon, J. M. Ryu, J.-P. Son, O. Seongil, H.-S. Yu, H. Lee, S. Y. Kim, Y. Cho, J. G. Kim, J. Choi, H.-S. Shin, J. Kim, B. Phuah, H. Kim, M. J. Song, A. Choi, D. Kim, S. Kim, E.-B. Kim, D. Wang, S. Kang, Y. Ro, S. Seo, J. Song, J. Youn, K. Sohn, and N. S. Kim, "25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications," in *International Solid-State Circuits Conference (ISSCC)*, 2021, pp. 350–352.
- [44] Y. Kwon, Y. Lee, and M. Rhu, "TensorDIMM: A Practical Near-Memory Processing Architecture for Embeddings and Tensor Operations in Deep Learning," in *International Symposium on Microarchitecture (MICRO)*, 2019, pp. 740–753.
- [45] J. H. Lee, H. Zhang, V. Lagrange, P. Krishnamoorthy, X. Zhao, and Y. S. Ki, "SmartSSD: FPGA Accelerated Near-Storage Data Analytics on SSD," *IEEE Computer Architecture Letters*, vol. 19, no. 2, pp. 110–113, 2020.
- [46] S. Lee, S.-h. Kang, J. Lee, H. Kim, E. Lee, S. Seo, H. Yoon, S. Lee, K. Lim, H. Shin, J. Kim, O. Seongil, A. Iyer, D. Wang, K. Sohn, and N. S. Kim, "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology : Industrial Product," in *International Symposium on Computer Architecture (ISCA)*, 2021, pp. 43–56.
- [47] V. T. Lee, A. Mazumdar, C. C. del Mundo, A. Alaghi, L. Ceze, and M. Oskin, "Application codesign of near-data processing for similarity search," in *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2018, pp. 896–907.

- [48] T. S. Lehman, A. D. Hilton, and B. C. Lee, "PoisonIvy: Safe speculation for secure memory," in *International Symposium on Microarchitecture (MICRO)*, 2016, pp. 1–13.
- [49] H. Li, T. F. Wu, S. Mitra, and H.-S. P. Wong, "Resistive RAM-centric computing: Design and modeling methodology," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 64, no. 9, pp. 2263–2273, 2017.
- [50] B. Libert, T. Peters, M. Joye, and M. Yung, "Linearly homomorphic structure-preserving signatures and their applications," *Designs, Codes and Cryptography*, vol. 77, no. 2, pp. 441–477, 2015.
- [51] J. Liu, H. Zhao, M. A. Ogleari, D. Li, and J. Zhao, "Processing-in-memory for energy-efficient neural network training: A heterogeneous approach," in *International Symposium on Microarchitecture (MICRO)*, 2018, pp. 655–668.
- [52] Q. Liu, B. Gao, P. Yao, D. Wu, J. Chen, Y. Pang, W. Zhang, Y. Liao, C.-X. Xue, W.-H. Chen, J. Tang, Y. Wang, M.-F. Chang, H. Qian, and H. Wu, "33.2 a fully integrated analog rram based 78.4tops/w compute-in-memory chip with fully parallel mac computing," in *International Solid-State Circuits Conference*, 2020, pp. 500–502.
- [53] A. López-Alt, E. Tromer, and V. Vaikuntanathan, "On-the-fly multiparty computation on the cloud via multikey fully homomorphic encryption," in *Symposium on Theory of computing (STOC)*, 2012, pp. 1219–1234.
- [54] D. A. McGrew and J. Viega, "The security and performance of the galois/counter mode (GCM) of operation," in *Progress in Cryptology (INDOCRYPT)*, 2005, pp. 343–355.
- [55] P. Mohassel and P. Rindal, "ABY3: A mixed protocol framework for machine learning," in *Conference on Computer and Communications Security (CCS)*, 2018, p. 35–52.
- [56] S. Na, S. Lee, Y. Kim, J. Park, and J. Huh, "Common counters: Compressed encryption counters for secure GPU memory," in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2021, pp. 1–13.
- [57] L. Nai, R. Hadidi, J. Sim, H. Kim, P. Kumar, and H. Kim, "GraphPIM: Enabling Instruction-Level PIM Offloading in Graph Computing Frameworks," in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2017, pp. 457–468.
- [58] M. Naumov, D. Mudigere, H.-J. M. Shi, J. Huang, N. Sundaraman, J. Park, X. Wang, U. Gupta, C.-J. Wu, A. G. Azzolini *et al.*, "Deep learning recommendation model for personalization and recommendation systems," *arXiv preprint arXiv:1906.00091*, 2019.
- [59] P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," in *International conference on the theory and applications of cryptographic techniques*. Springer, 1999, pp. 223–238.
- [60] B. Reagen, W.-S. Choi, Y. Ko, V. T. Lee, H.-H. S. Lee, G.-Y. Wei, and D. Brooks, "Cheetah: Optimizing and accelerating homomorphic encryption for private inference," in *International Symposium on High-Performance Computer Architecture (HPCA)*, 2021, pp. 26–39.
- [61] P. Rogaway, "Authenticated-encryption with associated-data," in *Conference on Computer and Communications Security (CCS)*, 2002, pp. 98–107.
- [62] B. Rogers, S. Chhabra, M. Prvulovic, and Y. Solihin, "Using address independent seed encryption and bonsai merkle trees to make secure processors os-and performance-friendly," in *International Symposium on Microarchitecture (MICRO)*, 2007, pp. 183–196.
- [63] G. Saileshwar, P. J. Nair, P. Ramrakhiani, W. Elsasser, and M. K. Qureshi, "Synergy: Rethinking secure-memory design for error-correcting memories," in *International Symposium on High Performance Computer Architecture*, 2018, pp. 454–465.
- [64] S. Seshadri, M. Gahagan, S. Bhaskaran, T. Bunker, A. De, Y. Jin, Y. Liu, and S. Swanson, "Willow: A User-programmable SSD," in *Symposium on Operating Systems Design and Implementation (OSDI)*, 2014, pp. 67–80.
- [65] S. T. Setty, R. McPherson, A. J. Blumberg, and M. Walfish, "Making argument systems for outsourced computation practical (sometimes)," in *NDSS*, vol. 1, no. 9, 2012, p. 17.
- [66] Y. S. Shao, B. Reagen, G.-Y. Wei, and D. Brooks, "Aladdin: A pre-rtl, power-performance accelerator simulator enabling large design space exploration of customized architectures," in *International Symposium on Computer Architecture (ISCA)*, 2014, pp. 97–108.
- [67] W. Shi, H.-H. S. Lee, M. Ghosh, C. Lu, and A. Boldyreva, "High efficiency counter mode security architecture via prediction and precomputation," in *International Symposium on Computer Architecture (ISCA)*, 2005, pp. 14–24.
- [68] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership Inference Attacks Against Machine Learning Models," in *2017 IEEE Symposium on Security and Privacy (SP)*, 2017, pp. 3–18.
- [69] G. Singh, L. Chelini, S. Corda, A. J. Awan, S. Stuijk, R. Jordans, H. Corporaal, and A.-J. Boonstra, "Near-memory computing: Past, present, and future," *Microprocessors and Microsystems*, vol. 71, p. 102868, 2019.
- [70] G. Singh, D. Diamantopoulos, C. Hagleitner, S. Stuijk, and H. Corporaal, "NARMADA: Near-memory horizontal diffusion accelerator for scalable stencil computations," in *2019 29th International Conference on Field Programmable Logic and Applications (FPL)*, 2019, pp. 263–269.
- [71] Student, "The probable error of a mean," *Biometrika*, pp. 1–25, 1908.
- [72] G. E. Suh, D. Clarke, B. Gasend, M. Van Dijk, and S. Devadas, "Efficient memory integrity verification and encryption for secure processors," in *International Symposium on Microarchitecture (MICRO)*, 2003, pp. 339–350.
- [73] G. E. Suh, D. Clarke, B. Gassend, M. van Dijk, and S. Devadas, "AEGIS: Architecture for tamper-evident and tamper-resistant processing," in *ACM International Conference on Supercomputing 25th Anniversary Volume*, 2003, p. 357–368.
- [74] F. Tramèr and D. Boneh, "Slalom: Fast, verifiable and private execution of neural networks in trusted hardware," in *International Conference on Learning Representations (ICLR)*, 2019.
- [75] S. Volos, K. Vaswani, and R. Bruno, "Graviton: Trusted execution environments on gpus," in *Symposium on Operating Systems Design and Implementation (OSDI)*, Oct. 2018, pp. 681–696.
- [76] M. Wilkening, U. Gupta, S. Hsia, C. Trippel, C.-J. Wu, D. Brooks, and G.-Y. Wei, "RecSSD: Near Data Processing for Solid State Drive Based Recommendation Inference," in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2021, pp. 717–729.
- [77] Wm. A. Wulf and Sally A. McKee, "Hitting the memory wall: implications of the obvious," in *ACM SIGARCH Computer Architecture News*, vol. 23, no. 1. ACM, 1995, pp. 20–24.
- [78] C. Yan, D. Engleider, M. Prvulovic, B. Rogers, and Y. Solihin, "Improving cost, performance, and security of memory encryption and authentication," in *International Symposium on Computer Architecture (ISCA)*, 2006, p. 179–190.