

CIPHERMATCH: Accelerating Homomorphic Encryption-Based String Matching via Memory-Efficient Data Packing and In-Flash Processing

Mayank Kabra[†] Rakesh Nadig[†] Harshita Gupta[†] Rahul Bera[†] Manos Frouzakis[†]
Vamanan Arulchelvan[†] Yu Liang[†] Haiyu Mao[‡] Mohammad Sadrosadati[†] Onur Mutlu[†]
ETH Zurich[†] King’s College London[‡]

Homomorphic encryption (HE) allows secure computation on encrypted data without revealing the original data, providing significant benefits for privacy-sensitive applications. Many cloud computing applications (e.g., DNA read mapping, biometric matching, web search) use exact string matching as a key operation. However, prior string matching algorithms that use homomorphic encryption are limited by high computational latency caused by the use of complex operations and data movement bottlenecks due to the large encrypted data size. In this work, we provide an efficient algorithm-hardware codesign to accelerate HE-based secure exact string matching. We propose CIPHERMATCH, which (i) reduces the increase in memory footprint after encryption using an optimized software-based data packing scheme, (ii) eliminates the use of costly homomorphic operations (e.g., multiplication and rotation), and (iii) reduces data movement by designing a new in-flash processing (IFP) architecture.

CIPHERMATCH improves the software-based data packing scheme of an existing HE scheme and performs secure string matching using only homomorphic addition. This packing method reduces the memory footprint after encryption and improves the performance of the algorithm. To reduce the data movement overhead, we design an IFP architecture to accelerate homomorphic addition by leveraging the array-level and bit-level parallelism of NAND-flash-based solid-state drives (SSDs). We demonstrate the benefits of CIPHERMATCH using two case studies: (1) Exact DNA string matching and (2) encrypted database search. Our pure software-based CIPHERMATCH implementation that uses our memory-efficient data packing scheme improves performance and reduces energy consumption by 42.9× and 17.6×, respectively, compared to the state-of-the-art software baseline. Integrating CIPHERMATCH with IFP improves performance and reduces energy consumption by 136.9× and 256.4×, respectively, compared to the software-based CIPHERMATCH implementation.

1. Introduction

Homomorphic Encryption (HE) [1–9] offers a solution to perform computation directly on encrypted data without the need for decryption [5, 10, 11]. This capability is essential for processing large volumes of security-sensitive data (e.g., DNA sequences [12–17], biometrics [18, 19], and web application databases [20–23]) in shared computing environments. Many applications perform exact string matching on security-critical data, such as DNA string matching [13–17], biometric signature matching [19], and privacy-preserving database searches [21–23]. However, to perform the string matching operation in conventional systems (without HE), the data must be decrypted before computation (or stored in an unencrypted

format). Since cloud servers can be shared among multiple users, sensitive user data can become vulnerable to security threats and leaks [24–26]. HE can significantly benefit privacy-sensitive applications [27–31] that require exact string matching [13–17, 19, 21–23] as the fundamental operation by directly operating on encrypted data without requiring decryption.

Unfortunately, homomorphic operations are typically $10^4 \times$ to $10^5 \times$ slower than their traditional unencrypted counterparts in existing systems [32]. Prior works propose two main approaches to perform secure string matching: (1) the Boolean approach (e.g., [17, 33]), and (2) the arithmetic approach (e.g., [27, 29, 34]). The Boolean approach [17, 33] packs individual bits into a polynomial, encrypts it, and uses homomorphic XNOR and AND operations to perform secure string matching on a search pattern of any size. In contrast, the arithmetic approach [27, 29, 34] packs multiple bits into a polynomial, encrypts it, and employs homomorphic multiplication and addition operations to compute the Hamming Distance (HD) [35] for approximate or exact secure string matching on search patterns of only specific sizes (see §2.2 for more detail).

Despite recent performance improvements in accelerating HE operations (e.g., [32, 36, 37]), we make two key observations from prior works on secure string matching [17, 27, 29, 33, 34] that highlight opportunities for significant performance improvement. First, prior works that employ the Boolean approach incur substantial latency overheads due to the need to traverse large homomorphically encrypted databases (e.g., [17, 38]). In contrast, the arithmetic approach achieves a lower memory footprint after encryption using data packing schemes, but relies on costly homomorphic operations, such as homomorphic multiplication and rotation (see §3.1 for more detail). However, the data packing schemes and algorithms used in the arithmetic approach can be further optimized by using simpler operations (e.g., homomorphic addition) to significantly improve the performance of secure exact string matching.

Second, a significant portion of the performance overhead in HE comes from the movement of large homomorphically encrypted databases, which can be up to $1000 \times$ larger than their unencrypted counterparts [32, 37, 39], across the memory hierarchy. Several near-data processing approaches [40, 41], such as (1) processing-near memory (PnM) (e.g., [42–48]), (2) processing-using memory (PuM) (e.g., [49–52]), (3) in-storage processing (ISP) (e.g., [53–59]) and (4) in-flash processing (IFP) (e.g., [60–62]) have been proposed to address the data movement bottleneck between storage and main memory system and the compute-centric processors or accelerators. While these solutions aim to utilize high parallelism and high memory bandwidth of the storage and main memory devices, prior secure string matching techniques [17, 27, 29, 33, 34] have limitations in

the data packing schemes, resulting in large memory footprint leading to inefficient use of the existing near-data processing architectures, and compute-centric systems (see §3.1).

Our goal in this work is to improve the performance of homomorphic encryption-based secure string matching by providing an efficient algorithm-hardware codesign. To this end, we propose CIPHERMATCH¹, whose **key idea** is to reduce (1) the increase in memory footprint after encryption by optimizing the existing data packing schemes used by the arithmetic approach, (2) computational latency (i.e., execution time) of secure exact string matching by eliminating costly homomorphic multiplication or rotation operations, and (3) data movement overhead by moving computation into the storage devices using a new in-flash processing design.

Key Mechanisms. To enable CIPHERMATCH, we employ four key mechanisms. First, we develop a **memory-efficient data packing scheme** that packs *multiple bits* of the database and search pattern into a polynomial before encryption, reducing the memory footprint of the encrypted database and encrypted search pattern. Second, we implement secure matching using *only* the homomorphic addition operation (an optimization enabled by our efficient data packing scheme), which reduces the execution time of secure string matching algorithm. Third, we use in-flash processing and implement homomorphic addition directly within NAND-flash chips using read and latch operations [60, 62], utilizing the inherent array-level and bit-level parallelism in flash chips. Fourth, we utilize this in-flash processing architecture to enable CIPHERMATCH in a modern SSD-based storage system (called *CM-IFP*).

Key Results. We demonstrate the benefits of CIPHERMATCH using two case studies: (1) exact DNA string matching and (2) encrypted database search. We compare our pure software-based CIPHERMATCH implementation (CM-SW) that uses our memory-efficient data packing scheme on a real CPU system (see §5) with a state-of-the-art software baseline [27] (see §2.2). Our evaluation shows that CM-SW improves performance and reduces energy consumption by 42.9× and 17.6×, respectively, over the state-of-the-art. To assess the performance of CIPHERMATCH with IFP (CM-IFP), we compare CM-IFP with (1) CM-SW and (2) CIPHERMATCH implementation using processing-using memory (CM-PuM) [49]. Our results demonstrate that CM-IFP improves performance over CM-SW and CM-PuM by 136.9× and 1.4×, and reduces energy consumption by 256.4× and 3.3×, respectively.

This work makes the following **key contributions**:

- We propose CIPHERMATCH, a new approach that improves the performance and energy efficiency of homomorphic encryption (HE) based secure exact string matching.
- CIPHERMATCH introduces a new memory-efficient data packing scheme that drastically reduces the large increase in memory footprint after encryption, leading to faster computation and reduced storage overhead for HE operations. The proposed data packing scheme improves the performance and efficiency of both compute-centric and near-data processing architectures by increasing parallelism and reducing data

movement overhead of secure string matching.

- CIPHERMATCH is the first work to exploit in-flash processing for HE operations, enabling efficient secure string matching directly on encrypted data within the SSD. This approach reduces the data movement bottleneck and leverages the array-level and bit-level parallelism of flash memory.
- Both pure software-based and in-flash processing implementations of CIPHERMATCH significantly improve performance and energy efficiency over state-of-the-art software and hardware approaches on real-world applications like DNA string matching and encrypted database search.

2. Background

2.1. Homomorphic Encryption

Homomorphic encryption (HE) [1–9] is an encryption paradigm that allows users to compute directly on encrypted data (ciphertext) data without decrypting it. Most HE schemes build upon the *learning with errors (LWE)* problem [63, 64]. In the context of HE schemes, encrypted plaintexts can be typically represented in vector format, leading to highly parallel computation [65–67]. In this work, we utilize the Brakerski-Fan-Vercauteren (BFV) [9] HE scheme.

Brakerski-Fan-Vercauteren (BFV) HE Scheme. The BFV scheme [9] is built on the Ring-LWE problem [64], which leverages polynomial rings to compute on encrypted data. The BFV scheme operates within a ring structure $R = \mathbb{Z}_q[X]/(X^n+1)$, where (1) n (ring dimension) denotes the maximum polynomial degree, (2) q is the coefficient bit-length of ciphertext, and (3) \mathbb{Z}_q denotes a set of integers modulo q . Arithmetic operations within this ring are performed modulo $X^n + 1$. The plaintext space R_t consists of polynomials ($t \geq 2$), where t is the coefficient bit length of plaintext and (X^n+1) , while the ciphertext space R_q utilizes a larger coefficient modulus ($q \gg t$). These parameters (n, q, t) are defined according to the security level of the homomorphic encryption scheme [68]. The BFV scheme leverages several algorithms (e.g., key generation, encoding, encryption, and homomorphic operations) [5] to achieve its functionalities.

Key Generation. A user generates a secret key $sk \in \mathbb{P}_n(\mathbb{Z}_q)$ and a public key pair $(pk_0, pk_1) \in \mathbb{P}_n^1(\mathbb{Z}_q)$ based on the required security parameters. $\mathbb{P}_n(R)$ denotes the set of polynomials of degree less than n with coefficients in ring R .

Encoding and Encryption. A message $m \in \mathbb{Z}_t$ is converted to a plaintext polynomial $M(x) \in \mathbb{P}_n(\mathbb{Z}_t)$ using an encoding scheme. The encryption function takes the public key pair (pk_0, pk_1) and a plaintext polynomial $M(x)$ and outputs a ciphertext $C = (C_0, C_1)$.

$$C = (C_0(x), C_1(x)) \in (\mathbb{Z}_q[x], \mathbb{Z}_q[x]) \quad (1)$$

$$C_0(x) = [pk_0(x) + e_0(x) + M(x)]_q \quad (2)$$

$$C_1(x) = [pk_1(x) + e_1(x)]_q \quad (3)$$

$e_0(x)$ and $e_1(x)$ are polynomials introduced to achieve the security properties of the scheme.

Homomorphic Addition (Hom-Add). The addition between two ciphertext $C^1 = (C_0^1, C_1^1)$ and $C^2 = (C_0^2, C_1^2)$ is performed

¹CIPHERMATCH is derived from a combination of ciphertext from homomorphic encryption and string matching, reflecting its core functionality in secure string matching on encrypted data.

coefficient-wise on the corresponding polynomials in the ciphertexts generated after encryption, as shown in Eq. (4):

$$C(x) = (C_0^1(x) + C_0^2(x), C_1^1(x) + C_1^2(x)) \quad (4)$$

2.2. Secure String Matching Algorithm

The conventional string matching operation is widely used in various applications and is typically implemented with bitwise XNOR and AND operations [69, 70]. Database systems [21–23, 71] commonly use exact string matching to search a key (i.e., a given string) in a database (i.e., many entries consisting of key-value pairs). Various bioinformatics applications [13–17, 72, 73] (e.g., read mapping and alignment) map a genome to a reference genome to identify genetic variations, using a combination of exact and approximate string matching. These applications can be used in a privacy-sensitive context and thus can have strict data security rules for user data. To ensure data privacy and security, these applications, likely need to perform computations on encrypted data.

HE offers significant benefits for secure string search: 1) Low communication complexity: HE requires only two rounds of data exchange between client and server during computation and incurs minimal data transfer overhead compared to other approaches, such as Yao’s garbled circuit [74, 75]. 2) Non-interactiveness: Unlike multi-party computation (MPC) [75, 76], HE does not require continuous interaction of users during computation, reducing the need for extensive online communication. 3) Universality: HE allows the implementation of any string matching algorithm without extensive data preprocessing, making the data accessible for other computational tasks, unlike other approaches (e.g., [77–79]). 4) No data leakage: HE is designed to hide all information about encrypted data, except for maximum data size, making it a secure option for string matching on encrypted data. While symmetric searchable encryption-based string matching schemes (e.g., [80, 81]) can be very efficient in terms of storage overhead and computation speed, their leakage profiles tend to be quite substantial [82–84]. This has led to various attacks that manage to recover parts of the plaintext or the queries [84]. Hence, in this work, our focus is to leverage HE for secure string matching. There are two key approaches to performing string matching using HE.

Boolean Approach encrypts the individual bits of the database and the query using the TFHE (Fast Fully Homomorphic Encryption over the Torus) encryption scheme [85, 86]. In this context, homomorphic XNOR and homomorphic AND operations are performed on the encrypted database with the encrypted query. This method requires traversing an encrypted database and performing homomorphic operations on individual encrypted bits of the database and query. Each bit is encrypted and converted to an encrypted polynomial. The Boolean approach effectively limits noise growth in homomorphic encryption, enabling an arbitrary number of computations [85]. As a result, it supports string matching for queries of any length, providing users with flexibility in designing string matching algorithms.

Arithmetic Approach encrypts the database and the query using SHE (Somewhat Homomorphic Encryption)-based

schemes [27, 29, 34, 87]. In this context, we encode both an input database and a query as integers and follow a three-step procedure. First, we partition the integer M , where M is at most ‘ k ’ bits, and convert into a binary vector (M_0, \dots, M_{k-1}) which is then expressed as a polynomial $P(M) = \sum_{i=0}^{k-1} M_i x^i$. Second, we represent $P(M)$ as our new message and encrypt the polynomial (see §2.1). Third, we calculate the Hamming Distance (HD) [35] between encrypted database and encrypted query using homomorphic multiplication and addition. This method avoids encrypting individual bits separately and instead packs multiple bits (i.e., performs data packing) within a single ciphertext, significantly reducing the overall data size. However, this data packing scheme restricts the query size encoded within the plaintext polynomial, as SHE permits only a finite number of computations on encrypted data [9]. Hence, string matching can only be performed on specific query lengths. We discuss the limitations of the state-of-the-art techniques using this data packing method in §3.1.

2.3. Basics of NAND Flash Architecture

Figure 1 shows the internal organization of a modern 3D NAND-flash based SSD.

NAND Flash Organization. Multiple floating-gate transistors are stacked serially, forming a NAND string [88, 89]. A NAND string is connected to a bitline (BL), and NAND strings connected to multiple BLs compose a sub-block ①. The control gates of all cells at the same vertical position within a sub-block are linked to a single wordline (WL), allowing for the concurrent operation of these cells. A NAND flash block ② consists of several (e.g., 4 or 8) sub-blocks, and thousands of blocks comprise a plane ③. A block includes hundreds to thousands of pages ④, each of which is 4–16 KiB in size. The blocks in a plane share all the BLs in that plane, which implies that thousands of NAND strings share a single BL. A NAND flash chip ⑤ contains multiple (e.g., 2 or 4) dies ⑥, and each die contains multiple (e.g., 2 or 4) planes. Multiple dies can operate independently but share the command/data buses (i.e., channel ⑦) in a time-interleaved manner.

NAND Flash Operations and Peripheral Circuitry. Modern NAND flash memory operates in single-level (SLC) [90], multi-level (MLC) [91], triple-level (TLC) [92], or quad-level (QLC) [93] modes, where each cell stores 1, 2, 3, or 4 bits, respectively. In modern NAND flash memory, during a read operation [94–99] an initial precharge voltage (V_{pre}) charges the bitline, and a read voltage (V_{read}) is supplied to the targeted wordline. If the cell has a high threshold voltage (V_{th}), the bitline remains in the logical 1 state; otherwise, the bitline is pulled down to logical 0. During this phase, the bitline value is latched into the sensing latch ⑧ and sent to the SSD controller via flash channels. Modern NAND flash memory has a sensing latch and multiple (e.g., 3 or 4) data latches ⑨ per plane [88], to buffer data read from multiple bits stored in a single cell [100]. We explain in detail the functionality of both data latches (D-latches) and sensing latch (S-latch) in §4.3.

Modern SSD Organization. Modern SSDs consist of four key components: 1) the SSD controller ①, which includes multiple embedded cores to manage the flash translation layer (FTL) ②; 2) per-channel hardware flash memory controllers ③

for request handling [97, 101]; 3) DRAM (e.g., 2GB LPDDR4 DRAM for a 2TB SSD) ④ [102] for caching logical-to-physical (L2P) mappings and storing metadata; and 4) the host interface layer (HIL) ⑤ for communication between the host system and SSD. The FTL is responsible for I/O scheduling, address translation, and garbage collection and uses DRAM to cache a portion of L2P mappings at sub-byte granularity (with the memory overhead of $\sim 0.1\%$ of the SSD capacity), while the HIL handles I/O requests through AHCI [103] or NVMe [104] protocols.

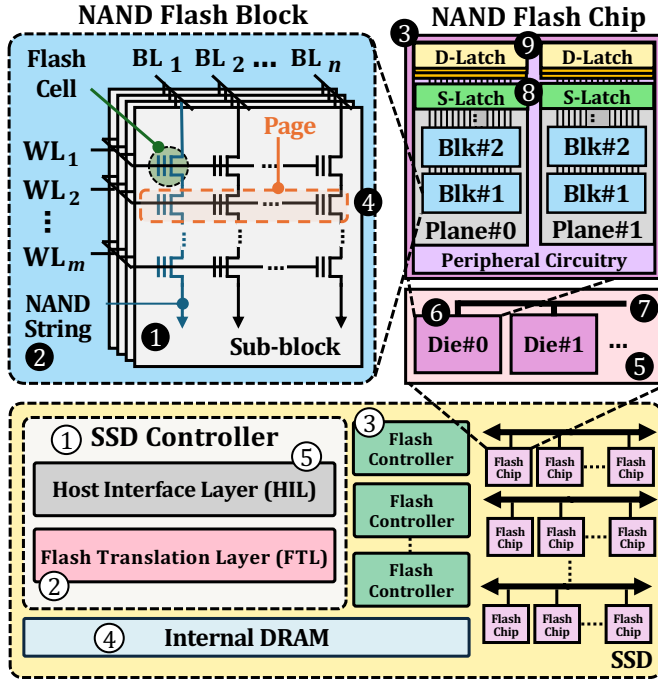


Figure 1: High-level overview of a modern SSD.

2.4. Near-Data Processing

Near-data processing (NDP) [40, 41] offers a promising solution to meet the high bandwidth requirements of memory-intensive workloads by performing computation where the data resides [105]. NDP can be categorized into two approaches: (1) processing-in (main) memory (PIM), which enables the memory subsystem to perform computation and (2) processing-in storage, which enables the storage devices to perform computation. PIM can be further categorized into two approaches: processing-using memory (PuM) and processing-near memory (PnM) [40]. Processing-in storage can be further categorized into two approaches: in-storage processing (ISP) and in-flash processing (IFP) [41].

Processing-using memory (PuM) approaches (e.g., [49–52]) use the operational principles of memory circuitry and chips to perform operations within the main memory, while processing-near memory (PnM) approaches (e.g., [42–48]) add processing elements inside or near the main memory to perform computations with reduced data movement overhead. However, with the increasing dataset sizes (e.g., [55, 106]), PIM approaches are increasingly bottlenecked by storage (SSD) accesses (as the entire dataset is typically stored in storage devices) [107]. With limited I/O bandwidth and high latency of SSD accesses, a

significant amount of time is spent loading the application data from storage to main memory [55, 60–62, 106].

In-storage processing (ISP) approaches (e.g., [53, 54, 56–59]) offload tasks to an in-storage accelerator or embedded cores in an SSD controller, which can significantly accelerate and improve the performance of data-intensive applications (e.g., [108–112]) by leveraging high internal SSD bandwidth to perform computation.

In-flash processing (IFP) approaches (e.g., [60–62, 113]) offer performance and energy benefits over ISP by exploiting the operational principles of NAND flash cells and circuitry for performing massively parallel computation, focusing on applications that use bulk-bitwise operations (e.g., [13, 49, 50, 114–131]).

Several NDP-based approaches have been proposed to accelerate homomorphic encryption (HE) by leveraging the massive parallelism and low-latency memory access of NDP architectures [39, 45, 132–136]. We discuss the limitations of these prior works in §3.3.

3. Motivation

We provide our motivational analysis in three areas. First, we qualitatively discuss the limitations of prior approaches [17, 27, 29, 33, 34] and quantitatively analyze two prior approaches (Boolean [17] and arithmetic [27]) that implement secure string matching using HE by comparing their memory footprint and execution time (see §3.1). Second, we provide an analysis to demonstrate the data movement overhead in the memory hierarchy for various database sizes (see §3.2). Third, we discuss the limitations of prior HE accelerators and near-data processing architectures for HE (see §3.3).

3.1. Limitations of Prior Approaches

Prior works typically leverage two key approaches to perform secure string matching (see §2.2): (i) Boolean and (ii) arithmetic. To compare these two approaches, we examine their four different characteristics: 1) *execution time*, the time required for the algorithm to perform the string matching operation; 2) *algorithm scalability*, the ability of the technique to scale with increasing database sizes; 3) *SIMD support*, whether the technique can leverage SIMD instructions to enhance performance; and 4) *flexible query size*, the ability of the technique to support different query sizes. Table 1 provides a comparison of four different characteristics of two different Boolean approaches [17, 33] with those of three different arithmetic approaches [27, 29, 34].

Boolean Approach. Prior works (e.g., [17, 33]) leverage the TFHE-rs library [86] to perform homomorphic Boolean operations (XNOR and AND) for secure string matching. Table 1 provides different characteristics of two state-of-the-art techniques: (1) Pradel et al. [33] use the TFHE library to develop a biometric string matching algorithm using basic Boolean homomorphic gates. However, their design does not leverage SIMD batching. (2) Aziz et al. [17] improve upon [33] by implementing SIMD batching to enhance the performance of Boolean homomorphic operations. However, two major challenges persist across Boolean approaches. First, performing a single secure exact string matching operation requires traversing the entire encrypted database and executing computationally expensive homomorphic operations, which significantly increase

the execution time. For instance, searching a 32-bit query in a small 32-byte database using HE on a real CPU system (see §5) takes 6.6s, whereas the same search on unencrypted data completes in only 5.9 μ s. Second, the memory footprint after encrypting individual bits (see §2.2) is larger (by more than 200 \times) compared to the memory footprint of unencrypted data.

Arithmetic Approach. Prior works (e.g., [27, 29, 34]) propose different data packing schemes to perform secure string matching using the arithmetic homomorphic operations (see §2.2). Table 1 provides different characteristics of three state-of-the-art techniques: (1) Yasuda et al. [27] utilize a data packing method that packs multiple bits in one ciphertext to achieve low latency for secure string matching operation. However, it suffers from two potential limitations. First, Hamming Distance (HD) [35] is used to find the string matches, which requires two costly homomorphic multiplications and three homomorphic additions for each comparison. Second, this approach divides the encrypted database into multiple ciphertexts and returns the same number of ciphertexts as a result of secure string matching. Hence, this approach is not scalable for large database sizes. (2) Kim et al. [34] address the scalability issue by employing a homomorphic equality (HomEQ) circuit (using homomorphic multiplication, rotation, and Frobenius map [137] operations) to return only the required results; however, it incurs high latency due to expensive HE operations. (3) Bonte et al. [29] improve upon [34] by implementing SIMD batching and compression techniques to enhance performance; however, it still requires performing the computationally expensive homomorphic operations and only supports specific query sizes.

Approach Type	Prior Work	Execution Time	Algorithm Scalability	SIMD Support	Flexible Query Size
Boolean Approach	[33]	High	✓	✗	✓
	[17]	High	✓	✓	✓
Arithmetic Approach	[27]	Low	✗	✗	✗
	[34]	High	✓	✗	✗
	[29]	High	✓	✓	✗

Table 1: Comparison of prior Boolean and arithmetic approaches based on (i) execution time, (ii) algorithm scalability (i.e., the ability of the technique to scale with increasing database size), (iii) SIMD support, and (iv) flexible query size (i.e., the ability of the technique to support different query sizes).

We quantitatively evaluate the memory footprint and execution time of two prior works, one following the Boolean approach [17] and another, the arithmetic approach [27]. Figure 2(a) shows the memory footprint of the encrypted database (Y-axis) for a given database size (X-axis) for these works. Figure 2(b) shows the execution time (in seconds) for secure string matching across different query sizes (X-axis) for the encrypted database (using the sizes from Figure 2(a)). Figure 2(c) shows the latency breakdown, highlighting the individual contributions of homomorphic addition and multiplication operations. We use a very small encrypted database size to understand the execution time of prior works without causing data movement in the memory hierarchy. We implement the Boolean approach [17] using the TFHE-rs library [86] and the arithmetic approach [27] using the Microsoft SEAL library [138] on a real CPU system (see §5). We make three key observations.

First, the increase in memory footprint introduced due to encryption is higher for the Boolean approach than the arithmetic approach. This is because the arithmetic approach efficiently groups multiple bits into a single polynomial using a data packing method before encryption (see §2.1). In contrast, the Boolean approach encrypts, each bit individually into a polynomial. Second, the execution time (in seconds) of the Boolean approach exceeds that of the arithmetic approach by 600 \times when performing secure string matching, averaged across various database sizes and query lengths. This is due to the higher overall cost of executing expensive Boolean homomorphic operations (e.g., XNOR and AND) on *individual encrypted bits*, compared to cheaper arithmetic homomorphic operations on *multiple packed bits*. Third, 98.2% of the latency in secure string matching using the arithmetic approach is due to the expensive homomorphic multiplication operations required to perform secure string matching.

Key Takeaway 1. *Data packing schemes used in arithmetic approaches offer significant performance benefits; however, they can be further optimized to minimize the use of costly homomorphic multiplication operations and to exploit highly parallel SIMD units to execute simpler homomorphic addition operations.*

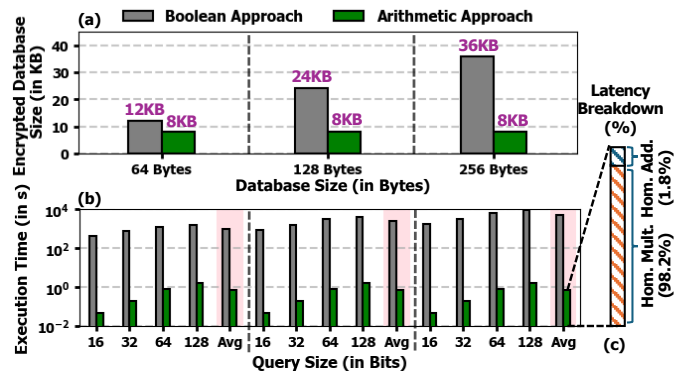


Figure 2: Comparison of one Boolean [17] and one arithmetic [27] technique in terms of (a) memory footprint, (b) execution time (in s), and (c) latency breakdown of the arithmetic approach [27].

3.2. Benefits of Near-Data Processing

Prior works (e.g., [55, 106]) demonstrate that the traditional string matching applications are constrained by data movement, as the time required for transferring large data volumes across the memory hierarchy significantly impacts performance. Once the databases are encrypted using HE, their memory footprint greatly increases (see §3.1), exacerbating the bottleneck caused by data movement. To understand the overheads of data movement from storage (SSD) to computation units, we profile the data movement overhead associated with transferring data to processing units for different database sizes, ranging from 8GB to 256GB. We assume three scenarios, where homomorphic secure string matching computations are performed in: (1) the CPU, (2) the main memory, and (3) the storage (SSD) controller.

Figure 3 shows the data transfer latency (normalized to the latency of transferring data to the CPU) from flash memory chips to (i) CPU, (ii) DRAM, and (iii) SSD controller. We observe that when the encrypted database is small (e.g., 8GB),

performing computations in main memory lowers data transfer latency by 25%, because doing so eliminates data movement to the CPU. As the database size increases, the overhead of transferring data from flash chips to DRAM becomes more significant, and thus diminishes the latency reduction benefits of performing computations in DRAM. For all database sizes, performing computations within the SSD controller reduces data transfer latency by over 80% compared to processing on the CPU. For the largest encrypted database size of 256GB, performing all computations in the SSD controller results in a 94% reduction in data transfer latency, whereas performing computations in DRAM provides only a 6% latency reduction.

Key Takeaway 2. *As encryption greatly increases the data size, applications become increasingly bound by data movement from storage. Performing computation closer to where the data resides, i.e., inside the SSD controller and flash memory can greatly reduce the data movement overheads.*

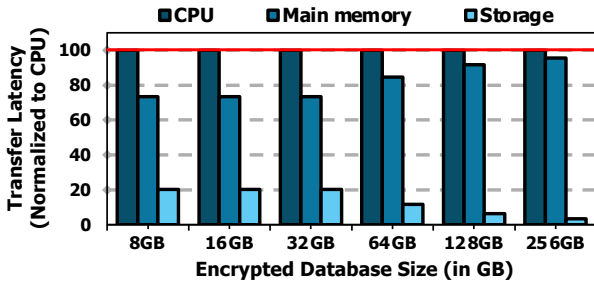


Figure 3: Transfer latency (normalized to the latency of transferring data to the CPU, denoted as $Y=100$) for loading data from flash chips to perform secure string matching in (1) CPU; (2) main memory; and (3) storage (SSD) controller. X-axis denotes the encrypted database size that is transferred.

Benefits of In-Flash Processing (IFP). IFP enhances performance and energy efficiency by utilizing the operational principles of flash arrays and cells for massively bit-parallel computation. Prior approaches (e.g., [60, 62]) use NAND flash memory logic to perform bulk-bitwise operations since these type of operations are widely used in many applications (e.g., [13, 49, 50, 114–131]) and are relatively simple to support in flash chips. Data movement is reduced by performing bulk bitwise operations in a massively parallel manner across all flash cells, arrays, and chips without moving the data outside flash chips, leading to high performance and large energy savings. Unfortunately, prior works do *not* investigate how to implement homomorphic operations using IFP. Although HE operations could potentially benefit from parallel processing using flash cells, arrays, and chips, they require support for arithmetic operations, which could be difficult to implement using IFP.

3.3. Limitations of Prior HE Accelerators

Existing HE accelerators adopt large on-chip caches or buffers (180MB for SHARP [139], 256MB for CraterLake [37], and 512MB for BTS [36]) to reduce frequent data loading from off-chip. However, such large on-chip storage may still be insufficient for performing string matching on large encrypted databases. Near-data processing in main memory or SSDs can be promising since it is challenging and costly to significantly increase either the off-chip memory bandwidth or the on-chip

storage. Prior works (e.g., [39, 45, 132–134]) propose frameworks to perform HE operations near main memory. However, data transfer from SSD to main memory can still limit the performance of secure string matching (see §2.4 & §3.2). Recent works [135, 136] accelerate arithmetic homomorphic operations by performing computations in the SSD controller. However, prior approaches do *not* optimize the data packing schemes, resulting in large memory footprints and inefficient utilization of the overall design (see §3.1). Moreover, computations in an SSD controller is still limited by the bandwidth of flash memory chips [60] and limited compute capability and parallelism due to power constraints in the SSD controller [55, 111, 128, 140].

3.4. Our Goal

We observe that there is a significant opportunity to exploit in-flash processing (IFP) to perform massively parallel homomorphic encryption-based secure exact string matching using bulk-bitwise computation capabilities of NAND flash memory chips. However, this requires efficient data packing techniques to enable bit-parallel computation and optimization strategies to reduce the complexity of homomorphic operations required for secure string matching.

Our goal is to develop an IFP-based hardware-software co-designed system that can perform efficient secure exact string matching. Specifically, we aim to (1) design a memory-efficient data packing scheme that enables bit-parallel operation and improves performance by eliminating the need for costly homomorphic multiplication operations (based on Key Takeaway 1), (2) exploit in flash processing to reduce the data movement bottleneck in secure homomorphic string matching (based on Key Takeaway 2), and (3) overcome the disadvantages of prior methods (Table 1 and §3.1).

4. CIPHERMATCH: Design

4.1. Overview

We propose **CIPHERMATCH**, an efficient algorithm-hardware co-designed system to accelerate secure string matching. **CIPHERMATCH algorithm** (see §4.2) comprises two key components: (1) a memory-efficient data packing scheme to reduce the memory footprint of encrypted data, and (2) a secure string matching algorithm that utilizes only homomorphic addition operations. **CIPHERMATCH hardware design** (see §4.3) has two key components: (1) a new in-flash processing (IFP) architecture to accelerate the secure string matching algorithm by reducing the data movement bottleneck, and (2) end-to-end system design to enable CIPHERMATCH in storage (SSD) devices.

4.2. Algorithm Design

CIPHERMATCH uses a memory-efficient data packing scheme to reduce memory footprint and implements a secure string matching algorithm with *only* homomorphic addition operations (i.e. avoiding costly homomorphic multiplications). For better understanding, we present the algorithm with specific parameters: $n = 1024$, $q = 32$, and $t = 16$ (see §2.1). With these parameters, the plaintext polynomial $P(M)$ has a degree of $n = 1024$, with each coefficient of size $t = 16$ bits. After encryption, the corresponding plaintext is converted to ciphertext polynomial $C = (C_0, C_1)$, where both C_0 and C_1 have a

degree of $n = 1024$, and each coefficient is $q = 32$ bits in size. Although we explain the algorithm with these specific parameters, the CIPHERMATCH algorithm can be adapted to any set of homomorphic encryption parameters defined by the HE standards [68].

4.2.1. Memory-Efficient Data Packing Scheme. Assume a binary string $P = (b_0, b_1, \dots, b_{k-1})$ of length k ; the string is partitioned into small segments, each containing t bits, where t is a predefined (HE) parameter (for our case, $t = 16$). If k exceeds t , the text is divided into multiple non-overlapping chunks of size t . For example, the first partition, $T^{(0)}$, will contain the first 16 bits: $T^{(0)} = (b_0, b_1, \dots, b_{15})$, and the second partition, $T^{(1)}$, will contain the next 16 bits: $T^{(1)} = (b_{16}, b_{17}, \dots, b_{31})$. We create a packed message $m(T)$ (see Eqn. (5)), which is simply the collection of these partitions:

$$m(T) = \left(T^{(0)}, T^{(1)}, \dots, T^{\left(\lfloor \frac{k}{t} \rfloor\right)} \right) \quad (5)$$

We convert the packed message $m(T)$ into a packed plaintext polynomial $M(x)$ with a maximum degree of n , where n is a predefined (HE) parameter (in our case, $n = 1024$). The polynomial representation is given by: $M(x) = \sum_{i=0}^{n-1} m_i x^i$, where m_i represents the packed bits at position i . If the number of elements in $m(T)$ exceeds n , we generate multiple plaintext polynomials. Specifically, if $m(T)$ contains L elements, we construct $\lceil L/n \rceil$ separate polynomials as shown in Eqn. (6):

$$M^{(j)}(x) = \sum_{i=0}^{n-1} m_{jn+i} x^i, \quad \text{for } j = 0, 1, \dots, \left\lceil \frac{L}{n} \right\rceil - 1. \quad (6)$$

We encrypt all the plaintext polynomials with a public key (pk) (see §2.1) to generate ciphertext $C^{(j)}(x)$ (see Eqn. (7)), with a maximum degree of n and coefficient size of q , where n, q are predefined (HE) parameters (in our case, $n = 1024, q = 32$):

$$C^{(j)}(x) = \text{Enc}(M^{(j)}(x), pk), \quad \text{for } j = 0, 1, \dots, \left\lceil \frac{L}{n} \right\rceil - 1. \quad (7)$$

Key Insight. Our memory-efficient data packing scheme reduces the memory overhead of encryption, resulting in the encrypted data being only $4\times$ larger (lower bound - assuming every coefficient of plaintext contains the maximum possible packed bits) than the original unencrypted data (as opposed to $64\times$ larger (lower bound) in the state-of-the-art data packing scheme [27]). Our technique achieves this $16\times$ memory overhead reduction by packing the maximum possible bits (in our case, 16 bits) into the plaintext coefficients, whereas prior work [27] uses a single-bit data packing approach. The $4\times$ increase in data size after encryption is due to a $2\times$ increase from converting plaintext into tuples and an additional $2\times$ increase in ciphertext coefficient size.

4.2.2. Secure String Matching Algorithm. We leverage the memory-efficient data packing scheme to design an efficient secure string matching algorithm using *only* homomorphic addition operations. Assume in plaintext domain, a binary query (Q); it is first negated ($\sim Q$) and added to a binary data (d) such that, if there is a match, the result will be a string of all 1's.

This mechanism can be adapted to the HE domain by replacing traditional addition operations with homomorphic addition operations to perform secure string matching. We utilize our memory-efficient data packing mechanism (§4.2.1) to efficiently pack the query and input data. Let the negated query ($\sim Q$) and input data (d) be represented as plaintext polynomials $M_Q(x)$ and $M_d(x)$, respectively, as shown in Eqn. (8):

$$M_{\sim Q}(x) = \sum_{i=0}^{n-1} \sim Q_i x^i, \quad M_d(x) = \sum_{i=0}^{n-1} d_i x^i \quad (8)$$

where $\sim Q_i$ and d_i are the respective packed binary coefficients (§4.2.1) of the query (Q) and the input (d). The polynomials are then encrypted, and the resulting ciphertexts are denoted as $C_{\sim Q}(x)$ and $C_d(x)$, respectively. After encryption, each coefficient of the ciphertext is mapped to a ring R_q , and the 16-bit coefficients are converted into 32-bit coefficients. While the coefficients are transformed during encryption, the original plaintext data stays in the same location; however, it is now represented by its encrypted value at that location. We perform homomorphic addition on the encrypted polynomials: $C_{\text{result}} = \text{Hom-Add}(C_{\sim Q}(x), C_d(x))$ where *Hom-Add* represents the homomorphic addition operation.

Index Generation. When a match occurs, the result of the homomorphic addition is a ciphertext where one (or more) coefficients is equal to the encrypted value corresponding to all 1's, i.e., $P_v(x) = 111\dots 11x^{1023} + \dots + 111\dots 11$. We call this polynomial $P_v(x)$ the "match polynomial". The encrypted result is compared with an encrypted "match polynomial" for each coefficient. This comparison checks whether the values in the encrypted data match those in the "match polynomial". If a match is found, it indicates that the query has successfully found a corresponding value in the input data. Based on this comparison, we determine the exact location where the match occurred.

We describe how to use this secure string matching algorithm within a client(user)-server model, as shown in Algorithm 1.

1) Database Preparation (lines 1-3). We assume that the database is initially flattened into a binary vector (lines 1-2), which is then transformed into an encrypted packed polynomial (line 3) using our memory-efficient data packing scheme before being stored on the server.

2) Query Preparation (lines 4-9). To efficiently perform secure string matching over a large encrypted database, we employ our memory-efficient data packing scheme that minimizes the memory footprint. We describe query preparation in four key steps that enable efficient parallel string matching.

- **Lines 4-6.** Given a client/user query (Q) of length 8 bits (e.g., 001...10) (line 4), we first negate it (line 5) and replicate it multiple times to construct a structured plaintext polynomial (line 6). This ensures that the same query is present across multiple coefficients of the same polynomial. As a result, when homomorphic addition is performed, multiple coefficient-wise additions perform multiple string matching operations (see §4.2.2) in parallel.
- **Line 7.** The replicated query is embedded within a structured polynomial $P_1(x)$, where each coefficient stores the same binary pattern: $P_1(x) = 110\dots 01110\dots 01x^{1023} +$

Algorithm 1 : CIPHERMATCH algorithm demonstrated in a client(user)-server model.

Database Preparation: (Server Side)

- 1: **Input:** Assume data $(d) = (10110001\dots 10)$ of length $m > n$.
 - 2: **Generate:** Packed vector $d = (d_0, \dots, d_{m'})$, where d_i are 16-bit partitions, $m' = \lceil m/16 \rceil$, and construct polynomials $P_j(d) = \sum_{i=j}^{j+1023} d_i x^i$ for $j \in (0, m'/1024)$.
 - 3: **Encrypt:** $E(P_j(d), \text{pk}) = C_j^s(x) = (ct_{0,j}^s, ct_{1,j}^s)$, store on the server.
-

Query Preparation: (Client/User Side)

- 4: **Input:** Assume query $(Q) = (001\dots 10)$ of length $y \leq m$.
 - 5: **Generate:** $Q' = \sim Q$ as a packed vector $pv = (Q_0, \dots, Q_{\frac{y}{16}})$, where Q_i represents 16 bits.
 - 6: If $\frac{y}{16} < 1024$, repeat the vector (pv, pv, \dots, pv) to fill the coefficients of the plaintext polynomial.
 - 7: Construct a polynomial $P(Q') = \sum_{i=0}^{1023} Q_i x^i$
 - 8: Perform $Q'_i < 1$ for n times and repeat steps 6-8.
 - 9: **Encrypt:** $E(P(Q'), \text{pk}) = C_i^u(x) = (ct_{0,i}^u, ct_{1,i}^u)$. Send this encrypted query to the server.
-

Secure String Search: (Server Side)

- 10: **Homomorphic Addition:** $E(P_i(Q') + P_j(d)) = (ct_{0,i,j}^s + ct_{0,i}^u, ct_{1,i,j}^s + ct_{1,i}^u)$ for all i and j .
 - 11: **Output:** This step generates multiple resultant ciphertexts $(ct_{0,i,j}^r, ct_{1,i,j}^r)$ between server (s) and user (u) ciphertexts.
 - 12: **Index Generation:** If the ciphertext contains a "match polynomial" in any coefficient, a match is found. The location index is generated and sent back to the client/user.
-

$110\dots 01110\dots 01x^{1022} + \dots + 110\dots 01110\dots 01$. Each coefficient fits within $t = 16$ -bits. If the query length exceeds 16 bits, the bits are distributed across multiple coefficients while ensuring efficient packing.

- o **Line 8.** To detect all possible alignments of the query within the encrypted database, we generate multiple left-shifted variants of Q . This results in a set of shifted polynomials $P_1(x), P_2(x), \dots, P_8(x)$, where each polynomial represents a different alignment of the query.
- o **Line 9.** Each shifted polynomial is then homomorphically encrypted using the same HE parameters, producing ciphertext polynomials $C_1^u(x), C_2^u(x), \dots, C_8^u(x)$, where u represents the user (i.e., client) query. These encrypted polynomials are then used to find a match against the encrypted database.

3) Secure String Search (lines 10-12). These encrypted query polynomials are then transferred to the server, which performs an efficient homomorphic addition (*Hom-Add*) (lines 10-11) with the encrypted database. The server generates an index (line 12) by comparing the resultant polynomial with the encrypted "match polynomial".

4.3. Hardware Design

CIPHERMATCH introduces a new in-flash processing architecture (see §4.3.1) to perform efficient bitwise operations and enable bit-serial addition operations, using which we can efficiently implement homomorphic addition. To leverage the benefits of this IFP architecture, we discuss the end-to-end system design (see §4.3.2) to integrate CIPHERMATCH in an SSD.

4.3.1. In-Flash Processing (IFP) Architecture. Our IFP architecture is inspired by two recent IFP works, ParaBit [62] and Flash-Cosmos [60]. ParaBit enables parallel bitwise operations (AND, OR) by controlling the data and sensing latch circuitry within the flash chip. Flash-Cosmos leverages the existing XOR circuit in NAND flash peripherals to perform bitwise XOR operations and introduces an Enhanced SLC Programming (ESP) method for reliable execution of these operations. Additionally, Flash-Cosmos supports bulk bitwise computation of AND and OR operations across data stored in NAND flash memory. However, there are limitations with both prior approaches. First, ParaBit inherently allows data transfer only from the S-latch to the D-latch and performs bitwise AND and OR operations by controlling one-sided data movement, which restricts the reuse of intermediate or previously computed results. Second, Flash-Cosmos only supports bulk-bitwise operations across the data present in the NAND flash memory. To perform string matching with a query, it is necessary to write the query string to NAND flash memory cells, which adds latency and accelerates the wearout of the flash memory.

Figure 4 illustrates the design of the sensing latch (S-latch) and data latch (D-latch) in NAND flash memory, which supports bitwise (AND, OR, and XOR) operations. In this design, we build on a prior work [141], which adds two transistors (M7 and M8 in Figure 4) in the existing NAND flash peripheral circuit. We utilize these transistors to enable *bi-directional data flow* (shown as blue arrows in Figure 4) by controlling the activation of transistors (M1-M8). This enables fine-grained control over bitwise AND and OR operations by explicitly managing the data transfer direction, ensuring that results are stored in the desired latch while allowing efficient reuse of intermediate or previously computed results. We use TLC NAND flash memory to leverage its *multiple data latches* (see §2.3) for intermediate data storage while operating it in SLC mode to ensure reliable computation of bitwise operations using the ESP method [60]. We explain (i) the data transfer from the S-latch to the D-latch and (ii) the execution of bitwise (AND, OR, and XOR) operations using the proposed NAND flash peripheral circuitry (which fundamentally differs from ParaBit [62]) as follows.

1) Data transfer from S-latch to D-latch occurs in three key steps. First, data is read from the flash cell to the S-latch **1** using the conventional flash read operation (see §2.3). Second, the M6 transistor is activated (by setting RST_D **2** to logical 1) to reset OUT_D to logical 0. Third, enabling the M5 transistor (by setting SET_D **3** to logical 1) allows OUT_D to be determined by OUT_S: if OUT_S is at logical 1, the M4 transistor is activated, which pulls OUT_D to GND and sets OUT_D to logical 1; otherwise, OUT_D remains at logical 0. This is equivalent to copying data from the S-latch to the D-latch.

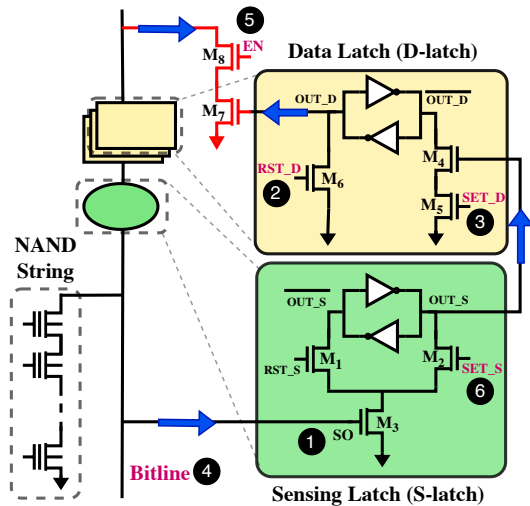


Figure 4: NAND flash peripheral circuitry with S-latch and D-latches to perform bitwise operations. Transistors (in red) show the modifications proposed by prior work [141]. The arrows (in blue) denote the data flow between the latches.

2) Bitwise AND of operands in S-latch and D-latch is performed in three key steps. First, we assume that the input data is present in the S-latch and the D-latch. Second, we precharge the bitline ④ and enable the M8 transistor (by setting EN ⑤ to logical 1). If the D-latch holds a logical 0, the M7 transistor does not turn on, keeping the bitline precharged at logical 1; otherwise, it pulls the bitline to GND, representing logical 0. Third, we enable M2 (by setting SET_S ⑥ to logical 1) and if the bitline value is logical 1, it pulls OUT_S to GND, making it logical 0, regardless of the previously stored value; otherwise, if the bitline value is logical 0, it turns off the M3 transistor, and OUT_S retains its previously stored value. This is equivalent to performing an AND operation between the values present in the D-latch and the S-latch and storing the result in *only* the S-latch.

3) Bitwise OR of operands in S-latch and D-latch is similar to performing data transfer from S-latch to D-latch. Instead of resetting the OUT_D value, the OUT_S value from the S-latch is transferred to the D-latch. We enable M5 (by setting SET_D ③ to logical 1) of the target D-latch. If OUT_S is logical 1, this pulls $\overline{\text{OUT_D}}$ to GND, setting OUT_D to 1. If OUT_S is logical 0, $\overline{\text{OUT_D}}$ and OUT_D retain their current values. This is equivalent to performing an OR operation between the values present in the D-latch and the S-latch and storing the result in *only* the D-latch.

4) Bitwise XOR of operands in D-latches is possible to perform using the existing peripheral circuit in most modern NAND flash chips [60, 142, 143]. An XOR circuit is present in between two D-latches to perform on-chip data randomization during write operations [144] or built-in error detection during chip testing [145]. In our design, we use this XOR circuitry to perform an XOR operation between D-latches 1 and 2, storing the result in D-latch 1. Flash-Cosmos [60] utilizes this XOR circuitry to perform Boolean operations.

CIPHERMATCH uses the proposed IFP architecture to enable the bit-serial addition mechanism to perform homomorphic addition across multiple bitlines in a flash array, processing data in parallel within each flash chip. We describe the required data

layout and the process for performing bit-serial addition.

Data Layout. To perform bit-serial addition for 32-bit operands (see §4.2.1), the key challenge is the carry propagation across all bit positions. In conventional NAND flash memory, operands are stored in a horizontal data layout, where the bits of each operand are placed contiguously in the cells of one wordline. When these bits are read, they are stored in the sensing latches connected to each bitline. This layout restricts the carry propagation because the carry bit must be moved between sensing latches as we perform the addition. To address this issue, we adopt a vertical data layout, where the bits of the operands are arranged along bitlines instead of wordlines [49]. In this layout, the bits of each operand are distributed across multiple wordlines, with each bit of a 32-bit element stored in a separate bitline. This enables us to compute the carry bit for each position while keeping it stored in a D-latch for the next bit position’s calculation, facilitating bit-serial addition [146].

Bit-Serial Addition. In the bit-serial addition process, a full-adder circuit is required to compute the sum and carry for each bit. The full-adder takes two operand bits, A_i and B_i , and an input carry, C_i , and produces a sum bit $S_i = A_i \oplus B_i \oplus C_i$ and a carry-out bit $C_o = (A_i \oplus C_i) \cdot B_i + A_i \cdot C_i$. The carry must then propagate through all bit positions as the addition proceeds. Figure 5 shows the sequence of operations for adding one bit of an input value B transferred to the SSD, to one bit of the value A stored in the flash memory. The carry-in value is initially stored in D-latch 2 and set to 0. ① Load the input bit B_i from the controller to the S-latch, and ② copy it to D-latch 1. ③ Perform the AND operation to calculate $B_i \cdot C_i$ in the S-latch and ④ perform the XOR operation to calculate $B_i \oplus C_i$ in D-latch 1. ⑤ Copy the result of the AND operation to D-latch 0. ⑥ Read the value A_i from the NAND flash cell and ⑦ copy it to D-latch 2. ⑧ Transfer the data from D-latch 1 to S-latch to perform $A_i \cdot (B_i \oplus C_i)$. ⑨ Perform the XOR operation between D-latch 1 and 2 to calculate $A_i \oplus B_i \oplus C_i$. ⑩ Copy $A_i \cdot (B_i \oplus C_i)$ to D-latch 2 and ⑪ copy $B_i \cdot C_i$ to the S-latch. ⑫ Perform the OR operation in D-latch 2 to calculate $(B_i \oplus C_i) \cdot A_i + B_i \cdot C_i$ as carry bit. ⑬ The sum bit S_i , which is stored in D-latch 1, is sent out to the SSD controller and the carry-out bit C_o is stored in D-latch 2. CIPHERMATCH calculates the next sum bit by repeating this procedure to calculate the sum of input B and input A.

Implementing Homomorphic Addition. In homomorphic encryption, ciphertexts (encrypted data) are represented as polynomials, with each coefficient stored as a large integer (in our case, 32 bits). Homomorphic addition performs an element-wise summation of these coefficients, which can be performed independently and in parallel. Our NAND flash architecture supports bit-serial addition, enabling efficient execution of homomorphic addition directly in flash memory. We organize the 32-bit coefficient of a polynomial along a single bitline and distribute polynomial coefficients across multiple bitlines. Bit-serial addition is performed within each bitline by serially summing the 32 bits of the stored coefficient with the coefficient of the query. To maximize throughput, we exploit multiple levels of parallelism: bitline-level parallelism across wordlines, chip-level parallelism across NAND chips, and channel-level parallelism for concurrent coefficient-wise addition operation.

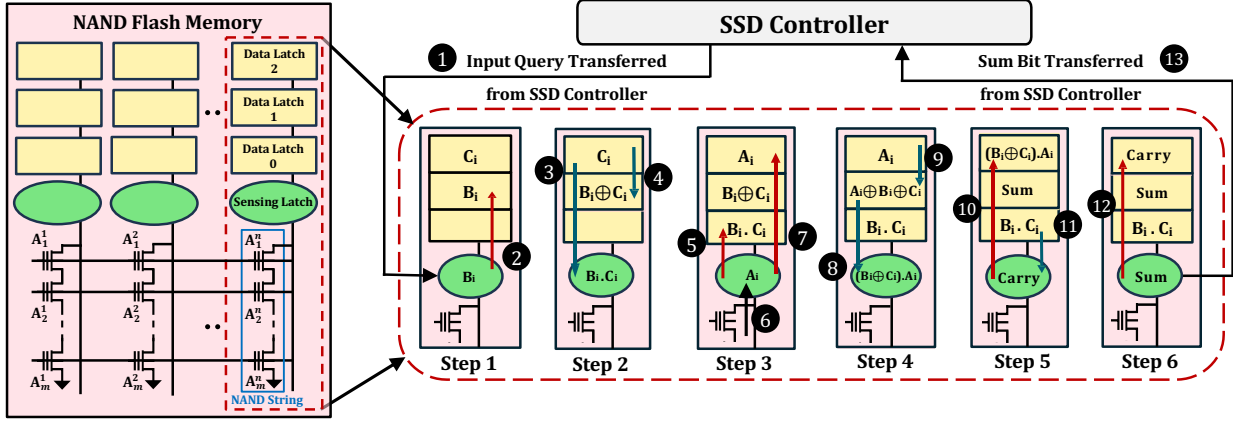


Figure 5: A single step of bit-serial addition using sensing and data latch circuitry in the flash array. Red, blue, and black lines denote data transfer from S-latch to D-latch, bitwise operations, and flash read operations, respectively.

Reliability. CIPHERMATCH maintains the reliability of bitwise operations and the underlying flash memory using two key strategies. First, CIPHERMATCH leverages enhanced SLC mode programming [60] to maximize the voltage difference between the two programming states within a cell, ensuring accurate data representation during bit-serial operations on individual bits. Second, CIPHERMATCH performs bit-serial addition completely using the latching circuit present in the flash chips, which avoids performing costly program/erase (P/E) cycle operations in the flash cell array that degrade the lifetime of flash memory.

4.3.2. End-to-End System Design. In this section, we present the end-to-end system design changes required to enable CIPHERMATCH in storage (i.e., in an SSD). We describe a system-level overview of the CIPHERMATCH algorithm, followed by a detailed explanation of the design of the individual components and system-level modifications to integrate CIPHERMATCH in an SSD.

System-level Overview. Figure 6 shows a system-level overview of CIPHERMATCH. We assume the database is first packed and encrypted with the optimized packing scheme (see §4.2.1) before being stored on the server. CIPHERMATCH follows a six-step procedure to perform secure string matching. ① The client(user) machine prepares the encrypted query, and the encrypted "match polynomial", and ② sends them to the server. ③ The server forwards them to the SSD controller (using our proposed system interface) and triggers a μ -program for homomorphic addition, consisting of a sequence of flash commands (see §4.3.1). ④ The μ -program executes homomorphic addition and utilizes array-level and bit-level parallelism of flash memory. ⑤ Finally, an index generation step (see §4.2.2) to identify the matching location is performed in the SSD controller and ⑥ the index is sent back to the client(user).

System Integration. We explain the necessary system changes to leverage CIPHERMATCH within the SSD from the host device.

1) Operations in Vertical Data Layout. The CIPHERMATCH algorithm relies on bit-serial addition, which requires the data to be stored vertically (we store the data across 32 wordlines to perform 32-bit addition). However, CPUs traditionally work with data stored in a horizontal layout. To address this challenge,

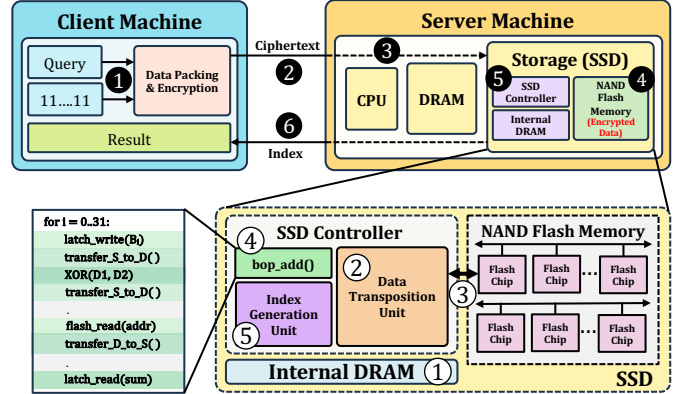


Figure 6: System-level overview of CIPHERMATCH.

CIPHERMATCH partitions the physical address space of an SSD into two regions: the conventional storage region and the CIPHERMATCH storage region. We handle the reads/writes for both regions separately. The conventional storage region handles read/write operations in TLC mode with a horizontal layout, while the CIPHERMATCH region operates in SLC mode with a vertical layout.

We explain the read and write operations in the CIPHERMATCH region as follows. For writes, CIPHERMATCH (i) encrypts the data, (ii) stores it in the DRAM buffer ①, (iii) performs data transposition in the SSD controller ② (see Figure 6) to prepare the data in a vertical layout, and (iv) write the vertical data to the flash chip ③. Each region maintains a separate mapping table for logical to physical address mapping, ensuring that data transposition remains transparent to the programmer. We introduce a software-based data transposition unit for this purpose (see Figure 6). The address mapping for the vertical layout is managed within its designated mapping table. Separate mapping tables handle commands for the conventional and CIPHERMATCH regions. For reads, the SSD controller reads NAND flash pages sequentially and stores them in the DRAM buffer. This data is transposed back into a horizontal layout using the software-based data transposition in the SSD controller before sending it back to the host. Prior works [108, 143] have explored such firmware modifications, and we provide an overhead analysis in §6.3.

2) Handling Page Faults to the CIPHERMATCH Region.

A page fault occurs when an application accesses data that is not present in DRAM, requiring the host to issue a read request to the SSD. In the CIPHERMATCH region, data is stored in the vertical layout, where each 32-bit coefficient is distributed across multiple bitlines (see §4.3.1). If the page fault happens to the CIPHERMATCH region, the SSD controller performs flash reads from multiple wordlines to read the homomorphically encrypted data and transposes the data to the conventional horizontal layout (using the data transposition unit). This read request experiences a longer latency due to reading *multiple* (e.g., 32) flash wordlines; however, we do not account for the latency for data transposition as it can be performed in parallel with flash reads (see §4.3.2). To handle these long-latency page faults efficiently, CIPHERMATCH leverages the operating system (OS) support for huge pages [147]. The page fault handler in the OS maintains a timeout (a configurable threshold) for these long latency read requests that defines the maximum wait time before retrying the request to prevent long delays.

3) Handling Dirty Writebacks to the CIPHERMATCH Region. A dirty writeback occurs when an application evicts modified data from the main memory (dirty data) that has not been written to the SSD. The operating system (OS) issues a write request to the SSD, where data is asynchronously written to the NAND flash memory. If the write request happens in the CIPHERMATCH region, the SSD controller transposes the data (using the data transposition unit) and writes it to NAND flash memory asynchronously at a granularity of 4KB (page size). CIPHERMATCH leverages the SSD’s existing write policies to manage these writes. Since these writes are asynchronous, data transposition has a minimal impact on the overall performance of the running application.

4) System Interfaces. We add three new commands to control the CIPHERMATCH region so that applications can leverage the benefits of CIPHERMATCH. We add a new flag to the conventional I/O read and write commands and design new interfaces, CM-read and CM-write, to handle data in the vertical format. When these new commands are used, the flag informs the SSD controller to transpose data while reading or writing the data. This 1-bit flag differentiates between conventional operations and the CIPHERMATCH region, activating the transposition unit as required. For search operations using CIPHERMATCH, we design a new command CM-search, which contains the encrypted query as a parameter. The host device transmits the CM-search command and encrypted query, which is then transferred to the data latches to initiate the bit serial addition (see §4.3.1).

Firmware Modifications. We explain the modifications to the SSD firmware to enable CIPHERMATCH operations.

1) μ -Program for Addition. CIPHERMATCH utilizes three key operations to perform bit-serial addition (see §4.3.1): (i) data transfer from S-latch to D-latch, (ii) Bitwise (XOR, AND and OR) operations between S and D latches, (iii) read operation to transfer the data from the flash cell to the S-latch. To efficiently enable the bit-serial addition operation, we introduce a new command, `bop_add` (bulk operation add) ④ in the flash translation layer (FTL). (see Figure 6). This command executes a sequence of operations that performs bit-serial addition across

multiple wordlines within NAND flash chips (see Figure 5), leveraging the inherent parallelism of NAND flash memory. CM-search command calls the `bop_add` command to perform the search operation.

2) Data Transposition Unit. CIPHERMATCH uses a software-based data transposition unit ② (see Figure 6) running on the SSD controller to handle CM-read and CM-write operations and page faults. The transposition unit operates at a 4KB granularity, corresponding to the NAND flash page size. We use a software-based data transposition unit instead of a hardware unit (see §7.1) because the latency of software-based data transposition latency using the SSD controller is $13.6\mu\text{s}$, which is lower than flash read latency [60, 148] and can be efficiently overlapped with flash read operations, minimizing the overall performance impact.

3) Index Generation Unit. The final index generation of the matched location is performed by the SSD controller ⑤ as part of the CM-search command. The result of the homomorphic addition is transferred from the D-latches to the SSD controller, which compares it with the encrypted “match polynomial” to generate the index of the matched location (see §4.2.2). We estimate the total latency of the index generation step running on the SSD controller cores (see §5) to be $3.42\mu\text{s}$, which can be effectively overlapped with the sequential flash read latency.

5. Methodology

Our evaluation methodology consists of three key components: (1) a real system evaluation methodology to assess the performance improvement provided by the software implementation of the CIPHERMATCH algorithm on an existing real CPU system (see §5.1), (2) a simulation-based evaluation methodology to analyze the benefits of implementing the CIPHERMATCH algorithm using IFP (see §5.2), and (3) details of evaluated workloads (see §5.3).

5.1. Real System Evaluation Methodology

Real System Configuration. Table 2 presents the configuration of the real CPU system used in our experiments. The system comprises of six out-of-order CPU cores, 32 GB of DDR4-2400 main memory, and a PCIe 4.0 NVMe SSD. We conduct our experiments to measure the increase in memory footprint and the performance of homomorphic operations.

CPU: Intel(R) Xeon(R) Gold 5118	<i>Microarchitecture:</i> Intel Skylake [149]
	x86-64 [150], 6 cores, out-of-order, 3.2 GHz
	<i>L1 Data + Inst. Private Cache:</i> 32kB, 8-way, 64B line
	<i>L2 Private Cache:</i> 256kB, 4-way, 64B line
	<i>L3 Shared Cache:</i> 8MB, 16-way, 64B line
Main Memory	32GB DDR4-2400, 4 channels
Storage (SSD)	Samsung 980 Pro PCIe 4.0 NVMe SSD 2 TB [102]
Operating System (OS)	Ubuntu 22.04.1 LTS

Table 2: Real CPU system configuration.

Evaluated Systems. We evaluate the software-based CIPHERMATCH algorithm (CM-SW) implemented using Microsoft SEAL [138] and compare its performance and energy consumption against two state-of-the-art approaches: (i) the Boolean approach [17] with TFHE-rs [86] and (ii) the arithmetic approach [27] with Microsoft SEAL [138].

Performance and Energy Measurements. We measure the execution time of our baselines and CM-SW using performance counters to focus specifically on the time spent executing homomorphic operations, excluding the overhead associated with encryption and encoding. To estimate energy consumption, we combine power parameters obtained from Intel’s RAPL tool [151, 152] with the execution time obtained from our performance measurements.

5.2. Simulation-Based Evaluation Methodology

Simulated Systems. We simulate four different systems, including CIPHERMATCH implemented on our IFP architecture (CM-IFP) using our in-house simulator to understand the trade-offs between compute-centric, memory-centric, and storage-centric processing approaches:

- **CM-SW (compute-centric)** is a conventional CPU implementation of the CIPHERMATCH algorithm. The goal is to identify the data movement bottlenecks and SIMD limitations in existing compute-centric systems.
- **CM-PuM (memory-centric)** is a processing-using memory approach in which we implement the CIPHERMATCH algorithm on an external 32GB DDR4 DRAM (see Table 3) and use the SIMDRAM framework [49] to perform string matching computations. The data is transferred from the SSD to the external DRAM, where parallel bit-serial addition operations are performed across multiple channels and banks. Using this bit-serial addition, homomorphic addition is executed to implement the CIPHERMATCH algorithm. The goal is to evaluate performance improvements by reducing data movement between the DRAM and the CPU while leveraging the parallelism of DRAM for computation.
- **CM-PuM-SSD (storage-centric)** is a processing-in storage approach in which we implement the CIPHERMATCH algorithm within the SSD’s 2GB LPDDR4 DRAM, following the SIMDRAM semantics (similar to CM-PuM). It utilizes the internal flash channel bandwidth for efficient data transfer between NAND flash memory and SSD-internal DRAM. The goal is to evaluate performance improvements by reducing data movement through external I/O while utilizing the parallelism of SSD-internal DRAM.

Performance Modeling. Table 3 presents different simulated near-data processing (NDP) system configurations. We develop an in-house simulator to model the SSD and external DRAM latency characteristics accurately. We model three different NDP systems and compare them against a simulated CM-SW² baseline.

- **For CM-PuM,** we develop a SIMDRAM-based simulator that models the latency of in-DRAM bitwise XOR, AND, and OR operations for implementing bit-serial addition. We use these values along with the DRAM configuration (see Table 3) in the simulator to compute the latency of bit-serial addition. Additionally, we model external I/O bandwidth and DRAM bandwidth to model data movement and its overheads accurately.

²For a fair comparison, we model the CPU-based system (CM-SW) by incorporating CPU computation, DRAM transfer, SSD read, and I/O transfer latency, based on the CPU configuration outlined in Table 2.

CM-PuM	32 GB DDR4-2400, 4 channel, 1 rank, 16 banks; Peak throughput: 19.2 GB/s
	Latency: T_{bbop} : 49 ns; Energy: E_{bbop} : 0.864 nJ; where <i>bbop</i> is bulk bitwise operation
	SSD External-Bandwidth: 7-GB/s external I/O bandwidth; (4-lane PCIe Gen4)
CM-IFP and CM-PuM-SSD	48-WL-layer 3D TLC NAND flash-based SSD; 2 TB
	SSD Internal DRAM: 2GB LPDDR4-1866 DRAM cache; 1 channel, 1 rank, 8 banks
	NAND-Flash Channel Bandwidth: 1.2-GB/s Channel IO rate
	Controller Cores: ARM Cortex-R5 series @1.5GHz; 5 Cores [153]
	NAND Config: 8 channels; 8 dies/channel; 2 planes/die; 2,048 blocks/plane; 196 (4×48) WLs/block; 4 KiB/page
	Latency: T_{read} (SLC mode): 22.5 μ s [60]; $T_{AND/OR}$: 20 ns [62]; $T_{latchtransfer}$: 20 ns [62]; T_{XOR} : 30 ns [60]; T_{DMA} : 3.3 μ s; T_{bit_add} (CM-IFP): 29.38 μ s
	Energy: E_{read} (SLC mode): 20.5 μ J/channel [60]; $E_{AND/OR}$: 10nJ/KB [62]; $E_{latchtransfer}$: 10nJ/KB [62]; E_{XOR} : 20nJ/KB [60]; E_{DMA} : 7.656 μ J/channel; E_{index_gen} (SSD controller): 0.18 μ J/page size; E_{bit_add} (CM-IFP): 32.22 μ J/channel

Table 3: Simulated system configurations.

- **For CM-PuM-SSD,** we compute the latency of bit-serial addition using the same SIMDRAM-based simulator. We use the SSD-internal DRAM configuration to compute the latency of bit-serial addition. Additionally, we model the NAND-flash channel bandwidth and SSD-internal DRAM bandwidth to model overall data movement and its overheads.
- **For CM-IFP,** we develop our custom simulator and model three key operations required to enable CIPHERMATCH (see §4.3.1): (i) bitwise XOR, AND, and OR operations, (ii) latch transfer latency, and (iii) NAND flash read latency (see Table 3). We input the NAND flash configuration to our simulator to calculate the latency of our bit-serial addition algorithm (see Table 3). For SSD controller tasks, such as data transposition and index generation (see §4.3.2), we develop a C program and run it in a QEMU environment [154] configured for the ARM Cortex-R5 series [153] to obtain the latency parameters. We calculate the latency of one-bit serial addition (T_{bit_add}) using Eqn. (9):

$$T_{bit_add} = T_{bop_add} + 2 * T_{DMA} \quad (9)$$

$$T_{bop_add} = T_{read} + 2 * T_{XOR} + 5 * T_{latchtransfer} + 4 * T_{AND/OR} \quad (10)$$

Energy Modeling. To model DRAM energy consumption, we use the power values based on a DDR4 power model [155, 156]. To model SSD energy consumption, we use the SSD power values of Samsung 980 Pro SSDs [102] and the NAND flash power values measured by Flash-Cosmos [60]. Using these power values, we estimate energy consumption of DRAM, SSD, and NAND flash using our performance models. To calculate the energy consumption of CM-IFP, we use an equation similar to Eqn. (10), replacing the time variable with the energy consumption of bulk operations. Additionally, we include the energy consumed by the SSD controller during the index generation phase (E_{index_gen}) (see Table 3). We calculate the energy consumption of one-bit serial addition (E_{bit_add}) using Eqn. (11):

$$E_{bit_add} = E_{bop_add} + 2 * E_{DMA} + E_{index_gen} \quad (11)$$

5.3. Workloads

We evaluate the performance of CIPHERMATCH using two real-world workloads: 1) exact DNA string matching and 2) encrypted database search.

1) *Exact DNA String Matching*: DNA sequence analysis relies on exact string matching in the seeding operation to extract smaller reads/bases from a reference genome [13–17, 72, 73]. The query sizes for these reads can range from 8 to 128 base pairs. For performance evaluation, we use a single query and vary the query sizes between 16 bits and 256 bits. We design a synthetic workload that uses a 32GB DNA database, which grows to 128GB (larger than the external DRAM size) once encrypted with our memory-efficient data packing scheme (see §4.2.1).

2) *Encrypted Database Search*: In scenarios where a client needs to search for specific records on a server while preserving privacy and security, secure string matching is performed on an encrypted database [21–23, 71]. For performance evaluation, we design a synthetic workload that varies the database size from 2GB to 32GB, which grows from 8GB to 128GB upon encryption with our optimized packing scheme. We simulate 1000 queries to assess how different database sizes affect performance under two conditions: when the entire encrypted database fits in DRAM and when it must be fetched repeatedly from the SSD.

6. Evaluation

We evaluate the effectiveness of CIPHERMATCH at improving the performance and energy consumption of secure string matching. §6.1 evaluates the performance and energy consumption of software-based CIPHERMATCH (CM-SW). §6.2 evaluates the performance and energy consumption of hardware-based CIPHERMATCH (CM-PuM, CM-PuM-SSD, and CM-IFP). §6.3 studies the overheads of implementing CM-IFP on off-the-shelf SSDs.

6.1. Software-Based CIPHERMATCH Analysis

6.1.1. Effect of Query Size. Figure 7 presents the speedup of the arithmetic approach [27] and CM-SW over the Boolean approach [17] (denoted as $Y=10^0$) for different query sizes ranging from 16 bits to 256 bits. For this analysis, we use a single query and a database size of 128GB (see §5.3).

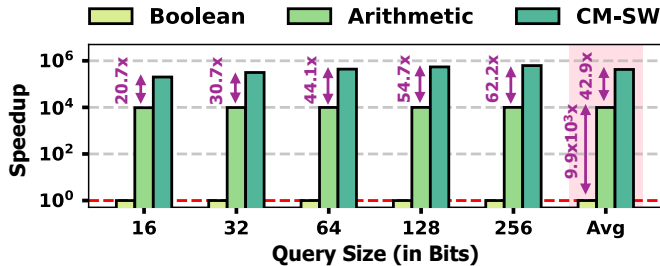


Figure 7: Speedup (higher is better) of the arithmetic approach [27] and CM-SW over the Boolean approach [17] (denoted as $Y=10^0$) for different query sizes, with a database size of 128GB and a single query.

We make three key observations. First, CM-SW outperforms the arithmetic [27] and Boolean [17] approaches by 20.7–62.2 \times and $2.0 \times 10^5 - 6.2 \times 10^5 \times$ for different query sizes, respectively.

Second, CM-SW speedup over the arithmetic approach [27] increases with query size. For example, CM-SW speedup increases from 20.7 \times to 30.7 \times when we increase the query size from 16 bits to 32 bits. This improvement is because CM-SW uses *only* homomorphic addition for secure string matching, which is computationally much less expensive than the arithmetic approach [27], which uses costly arithmetic homomorphic operations (e.g., multiplication) and more homomorphic operations are needed with an increase in query size. Third, our evaluation shows that CM-SW consistently outperforms prior works across all query sizes.

Figure 8 presents the energy consumption of the arithmetic approach [27] and CM-SW normalized to the Boolean approach [17] for different query sizes ranging from 16 bits to 256 bits. For this analysis, we use a single query and a database size of 128GB (§5.3).

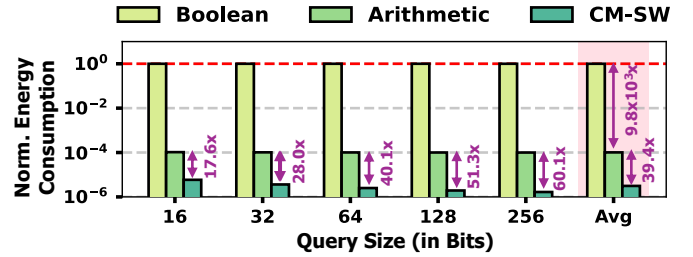


Figure 8: Energy consumption (lower is better) of the arithmetic approach [27] and CM-SW normalized to the Boolean approach [17] (denoted as $Y=10^0$) for different query sizes, with a database size of 128GB and a single query.

We make three key observations. First, CM-SW reduces the energy consumption by 17.6–60.1 \times and $1.6 \times 10^5 - 6.0 \times 10^5 \times$ compared to the arithmetic [27] and Boolean [17] approaches for different query sizes, respectively. Second, the energy consumption benefits of CM-SW over the arithmetic approach [27] increase with query size. For example, CM-SW reduces energy consumption from 17.6 \times to 28.0 \times over the arithmetic approach [27] when we increase the query size from 16 bits to 32 bits. Third, our evaluation shows that CM-SW consistently provides energy savings over prior works across all query sizes.

6.1.2. Effect of Database Size. Figure 9 presents the speedup of the arithmetic approach [27] and CM-SW over the Boolean approach [17] for different encrypted database sizes ranging from 8GB to 128GB. For this analysis, we use 1000 queries with a query size of 16 bits (see §5.3).

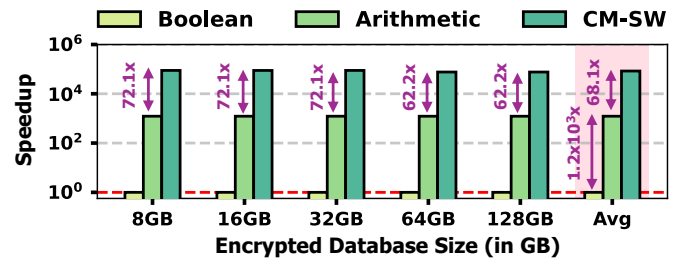


Figure 9: Speedup (higher is better) of the arithmetic approach [27] and CM-SW over the Boolean approach [17] (denoted as $Y=10^0$) for different encrypted database sizes, with a query size of 16 bits and 1000 queries.

We make three key observations. First, CM-SW outperforms the arithmetic [27] and Boolean [17] approaches by $62.2\text{-}72.1\times$ and $7.6\times 10^4 - 8.8\times 10^4\times$ for different encrypted database sizes, respectively. Second, when the database size exceeds 32GB, the performance improvement of CM-SW reduces by $1.16\times$. This is because the data movement overhead in CM-SW increases due to frequent I/O transfers for each query. In contrast, both the arithmetic [27] and Boolean [17] approaches saturate the memory footprint beyond the external DRAM size, even for smaller database sizes (e.g., 8 GB), due to their different data packing schemes. Third, our evaluation shows that CM-SW consistently outperforms prior works across all database sizes. **Summary.** We conclude that CM-SW provides significant performance improvement and energy savings over prior state-of-the-art approaches [17,27] across various query and database sizes. This improvement is primarily due to two reasons: (i) the use of *only* the homomorphic addition operation that is computationally less expensive than homomorphic multiplication and (ii) the use of a memory-efficient data packing scheme that reduces both memory footprint and data movement.

6.2. Hardware-Based CIPHERMATCH Analysis

6.2.1. Effect of Query Size. Figure 10 presents the speedup of CM-PuM, CM-PuM-SSD, and CM-IFP over CM-SW (denoted as $Y=10^0$) for different query sizes ranging from 16 bits to 256 bits. For this analysis, we use a single query and a database size of 128GB (§5.3).

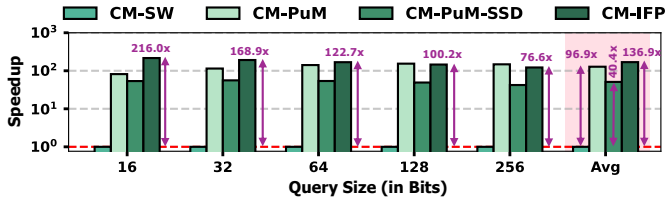


Figure 10: Speedup (higher is better) of CM-PuM, CM-PuM-SSD, and CM-IFP over CM-SW (denoted as $Y=10^0$) for different query sizes, with a database size of 128GB and a single query.

We make five key observations. First, CM-IFP, CM-PuM-SSD, and CM-PuM outperform CM-SW by $76.6\text{-}216.0\times$, $81.7\text{-}105.8\times$, and $26.4\text{-}53.9\times$, for different query sizes, respectively. These improvements are due to two reasons: (i) reduced data movement to the compute units and (ii) improved parallelism at the array-level and bit-level to perform homomorphic additions. Second, CM-IFP outperforms CM-PuM-SSD by $2.89\text{-}4.03\times$ for different query sizes. This improvement is due to two reasons: (i) data is not transferred outside the flash memory chips, reducing the data movement latency, and (ii) CM-IFP utilizes flash array, chip, and channel-level parallelism, which is higher than that of the SSD-internal DRAM to achieve higher throughput for computations.

Third, for small query sizes (e.g., 16-bit), CM-IFP outperforms CM-PuM by $2.64\times$; however, for large query sizes (e.g., 256-bit), CM-PuM achieves $1.21\times$ higher performance than CM-IFP. We observe that CM-IFP speedup over CM-SW varies with query size. For example, CM-IFP speedup reduces from $216.0\times$ to $168.9\times$ over CM-SW when we increase the query size from 16 bits to 32 bits. This is because for one query, even after reducing the data movement, performing computations for

a large query size (e.g., 256 bits) using DRAM reads are significantly faster than flash reads. Fourth, CM-PuM outperforms CM-PuM-SSD by $1.5\text{-}3.5\times$ for different query sizes. While CM-PuM-SSD reduces data movement by performing computations within the SSD and utilizing NAND-flash channel bandwidth, its performance is limited by the smaller internal DRAM, restricting the throughput of homomorphic additions compared to external DRAM. Fifth, our evaluation shows that CM-IFP provides the highest average performance across all query sizes.

Figure 11 presents the energy consumption of CM-PuM, CM-PuM-SSD, and CM-IFP normalized to CM-SW for different query sizes ranging from 16 bits to 256 bits. For this analysis, we use a single query and a database size of 128GB (see §5.3).

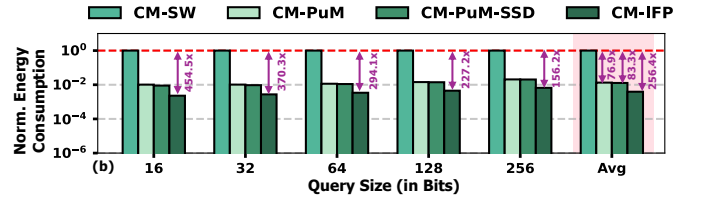


Figure 11: Energy consumption (lower is better) of CM-PuM, CM-PuM-SSD, and CM-IFP normalized to CM-SW (denoted as $Y=10^0$) for different query sizes, with a database size of 128GB and a single query.

We make three key observations. First, CM-IFP, CM-PuM-SSD, and CM-PuM reduce the energy consumption compared to CM-SW by $156.2\text{-}454.5\times$, $49.1\text{-}111.8\times$, and $48.6\text{-}98.3\times$ for different query sizes, respectively. Second, while CM-PuM outperforms CM-PuM-SSD in terms of performance (see Figure 10), CM-PuM-SSD is $1.06\times$ more energy efficient on average than CM-PuM. This is due to the lower energy cost of data transfers through internal NAND flash channels compared to external I/O channels. Third, our evaluation shows that CM-IFP provides the highest average energy savings across all query sizes.

6.2.2. Effect of Database Size. Figure 12 presents the speedup of CM-PuM, CM-PuM-SSD, and CM-IFP over CM-SW for different encrypted database sizes ranging from 8GB to 128GB. For this analysis, we use 1000 queries with a query size of 16 bits (see §5.3).

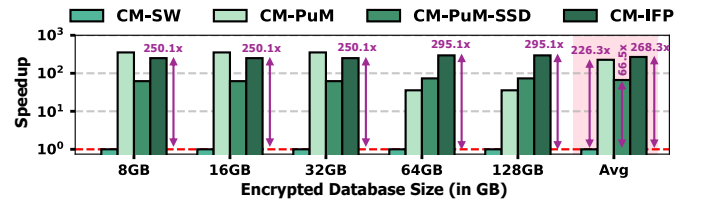


Figure 12: Speedup (higher is better) of CM-PuM, CM-PuM-SSD, and CM-IFP over CM-SW (denoted as $Y=10^0$) for different encrypted database sizes, with a query size of 16 bits and 1000 queries.

We make four key observations. First, CM-IFP, CM-PuM-SSD, and CM-PuM outperform CM-SW by $250.1\text{-}295.1\times$, $52.8\text{-}62.3\times$, and $35.6\text{-}353.4\times$ for different query sizes, respectively. Second, we compare CM-IFP and CM-PuM to understand the impact of SSD-external bandwidth. For small encrypted databases ($\leq 32\text{GB}$, which fits in external DRAM), CM-PuM

outperforms CM-IFP by $1.41\times$; however, for large encrypted databases ($> 32\text{GB}$, exceeding the size of external DRAM), CM-IFP achieves $8.29\times$ higher performance than CM-PuM. This is because, for smaller database sizes, the data can be loaded and then mostly accessed from DRAM for multiple queries and CM-PuM amortizes the cost of external I/O transfers to the SSD. However, for a large database that exceeds the DRAM capacity, the data must be fetched from the SSD, leading to frequent data movement between DRAM and SSD that degrades performance with CM-PuM.

Third, we compare CM-PuM and CM-PuM-SSD to understand the impact of NAND-flash channel bandwidth. For small encrypted databases ($\leq 32\text{GB}$, which fits in external DRAM), CM-PuM outperforms CM-PuM-SSD by $6.6\times$; however, for large encrypted databases ($> 32\text{GB}$, exceeding the size of external DRAM), CM-PuM-SSD achieves $1.75\times$ higher performance than CM-PuM. This is because, for smaller database sizes, the limited capacity of internal DRAM restricts the throughput of CM-PuM-SSD in homomorphic addition operations. In contrast, CM-PuM leverages extensive channel and bank-level parallelism. However, for a large database and multiple queries, CM-PuM-SSD utilizes NAND-flash channel bandwidth, reducing data movement time and outperforming CM-PuM, which requires frequent external I/O transfers to the SSD. Fourth, our evaluation shows that CM-IFP provides the highest average performance across all database sizes.

Summary. We draw two major conclusions from our analysis of hardware-based CIPHERMATCH implementations. First, all three hardware-accelerated CIPHERMATCH implementations (CM-PuM, CM-PuM-SSD, and CM-IFP) provide significant performance improvements and energy savings over the pure software-based CIPHERMATCH implementation (CM-SW) across various query and database sizes. This improvement is primarily due to (i) reduced data movement from/to main memory and/or storage devices and (ii) a significant increase in throughput by utilizing array-level and bit-level parallelism to perform homomorphic addition operations. Second, the IFP implementation of CIPHERMATCH (CM-IFP) provides the highest performance improvements and energy savings across all query and database sizes.

6.3. Overhead Analysis of CM-IFP

Storage Overhead. CM-IFP incurs three types of storage overhead. First, CIPHERMATCH partitions the SSD into two regions, i.e., SLC mode for the CIPHERMATCH region and TLC mode for the conventional region (see §4.3.2), reducing the overall capacity of the SSD. Second, the IFP bit serial algorithm sends the results to the controller for index generation (see §4.3.2), which requires additional storage in the SSD-internal DRAM. Based on our SSD configuration (see Table 3), CIPHERMATCH occupies 0.5 MB of SSD-internal DRAM space (4KB (page size) $\times 8$ (number of channels) $\times 8$ (number of dies) $\times 2$ (number of planes)) for storing the results of homomorphic addition. Third, CIPHERMATCH stores the μ -program (*bop_add*) in the internal-SSD DRAM, which requires less than 1KB .

Area Overhead. CM-IFP can be performed with minimal hardware modification to the NAND flash memory. CIPHERMATCH is built on two key designs: 1) hardware modifications

proposed by ParaBit [62] to enable bitwise (AND and OR) operations, which incurs 0.6% area overhead to modify the NAND-flash peripheral circuitry, and 2) Flash-Cosmos’s ESP technique [60] (see §4.3.1) to enable reliable bitwise operations, which does *not* incur any area overhead. Hence, the area overhead is $\sim 0.6\%$ of the die area of the baseline NAND flash memory.

7. Discussion

7.1. Hardware-Based Data Transposition Unit

CIPHERMATCH uses a software-based data transposition unit to perform data transposition using the SSD’s general-purpose core. Data transposition latency can be effectively overlapped with flash read latency ($22.5\mu\text{s}$) by pipelining the transposition operation with multiple flash reads. However, with advancements in flash memory technology that further reduces the read latency (e.g., $3\mu\text{s}$ in the latest Z-NAND flash memory [148]), we need to accelerate data transposition to hide the transposition latency during flash reads. To this end, we propose the integration of a dedicated hardware transposition unit adjacent to the SSD controller. The design of this unit is similar to SIMDRAM’s data transposition unit [49]. We implement the hardware-based data transposition unit in Verilog HDL and synthesize our design using a 22nm CMOS technology node [157] to estimate the latency and area overheads. The data transposition unit performs the transposition of 4KB data in 158ns and incurs an area overhead of 0.24 mm^2 .

7.2. Mitigation Techniques for Privacy Concerns

CIPHERMATCH returns the index of the matched location to the user after performing secure string matching (see §4.2.2), and this index must be securely transmitted from the server to the user to ensure privacy. To securely transfer the index of the matched location, we leverage a hardware-based 256-bit AES module present in commodity SSDs [158, 159] to encrypt the index before transmission across vulnerable channels (e.g., from an SSD in the cloud to the client system). During an offline step, the SSD controller generates a new AES key and sends the key to the client system. To secure the AES key during transfer, the SSD controller encrypts it with a public-key encryption algorithm before sending it to the client, mitigating the risk of exposure over vulnerable channels. This encryption occurs in an offline step, allowing the cost of key transfer to be amortized. We encrypt the index with the AES key and send the encrypted index back to the user/client. We implement the AES hardware unit (which works at 16 bytes granularity) using Verilog HDL and synthesize it using 22nm CMOS technology node [157]. The latency to encrypt 16 bytes is estimated at 12.6ns , with an area overhead of 0.13 mm^2 .

8. Related Work

To our knowledge, CIPHERMATCH is the first in-flash processing system designed to accelerate secure exact string matching. We have already qualitatively and quantitatively compared both software and hardware implementations of CIPHERMATCH to the state-of-the-art Boolean [17, 33] and arithmetic [27, 29, 34] approaches (see §3.1 & §6). In this section, we discuss other related works.

Secure String Matching. Prior works (e.g., [80, 160–164]) have explored symmetric searchable encryption (SSE) [82] for secure string matching, leveraging conventional encryption primitives [165, 166]. Faber et al. [160] improve SSE-based secure string matching by using an optimized data structure like partition trees [167] to represent the database and use conventional encryption techniques to perform secure string matching. Hahn et al. [161] enable efficient and secure substring searches over encrypted data outsourced to untrusted servers by using frequency-hiding order-preserving encryption (OPE) [168] techniques. Mainardi et al. [162] improve the secure string matching algorithm to achieve a sub-linear (polylogarithmic) data transfer (communication) complexity between the client and the server. Unfortunately, SSE-based schemes exhibit substantial leakage profiles, making them vulnerable to plaintext recovery attacks [84]. In contrast, homomorphic encryption (HE) [1] is based on strong mathematical foundations and has been proven to be secure [68]. CIPHERMATCH leverages homomorphic encryption for secure string matching and develops new software and hardware techniques to reduce HE’s computational and memory overheads.

Near-Data Processing for HE. Prior works on memory-centric computing (e.g., [39, 45, 132, 133, 169]) focus on performing homomorphic operations within main memory to reduce data movement overheads. However, these works can still be limited by data movement from/to storage devices, especially when the dataset size exceeds the size of main memory. Other works utilize FPGAs (e.g., [135, 136]), to perform homomorphic operations near storage, reducing the data movement from storage devices to the CPU. These solutions are constrained by limited computational capability and NAND-flash channel bandwidth, resulting in high latency and energy consumption for large datasets. CIPHERMATCH addresses these challenges by performing homomorphic addition operations *directly within* the NAND flash memory using operational principles of flash memory circuitry, thereby both reducing data movement and efficiently exploiting large-scale array-level and bit-level parallelism.

9. Conclusion

We introduce CIPHERMATCH, an algorithm-hardware co-designed system for accelerating secure exact string matching in solid-state drives (SSDs). CIPHERMATCH consists of two key components: (1) a memory-efficient data packing scheme that reduces the memory footprint after encryption and eliminates the use of costly homomorphic operations (i.e., multiplication), and (2) an in-flash processing (IFP) architecture that exploits array-level and bit-level parallelism in NAND flash memory to enable parallel string matching through homomorphic computations, requiring only small modifications to NAND flash chips in commodity SSDs. Our evaluation shows that both pure software-based and hardware-accelerated CIPHERMATCH implementations provide large performance improvements and energy savings over state-of-the-art approaches for homomorphic encryption-based secure exact string matching.

Acknowledgments

We thank the anonymous reviewers of ASPLOS 2025 for feedback. We thank the SAFARI group members for feedback

and the stimulating intellectual environment they provide. We acknowledge the generous gifts from our industrial partners, including Google, Huawei, Intel, and Microsoft. This work is supported in part by the ETH Future Computing Laboratory (EFCL), Huawei ZRC Storage Team, Semiconductor Research Corporation, AI Chip Center for Emerging Smart Systems (ACCESS), sponsored by InnoHK funding, Hong Kong SAR, and European Union’s Horizon programme for research and innovation [101047160 - BioPIM].

References

- [1] M. Ogburn, C. Turner, and P. Dahal, “Homomorphic Encryption,” *PROCS*, 2013.
- [2] D. Tourky, M. ElKawagy, and A. Keshk, “Homomorphic Encryption The “Holy Grail” of Cryptography,” in *ICCC 2016*.
- [3] C. Gentry and S. Halevi, “Implementing Gentry’s Fully Homomorphic Encryption Scheme,” in *EUROCRYPT 2011*.
- [4] A. Al Badawi, C. Jin, J. Lin, C. F. Mun, S. J. Jie, B. H. M. Tan, X. Nan, K. M. M. Aung, and V. R. Chandrasekhar, “Towards The Alexnet Moment For Homomorphic Encryption: HCNN, The First Homomorphic CNN on Encrypted Data With GPUs,” *TETC 2020*.
- [5] C. Gentry, “Fully Homomorphic Encryption using Ideal Lattices,” in *STOC 2009*.
- [6] M. Van Dijk, C. Gentry, S. Halevi, and V. Vaikuntanathan, “Fully Homomorphic Encryption over the Integers,” in *EUROCRYPT 2010*.
- [7] D. Boneh, C. Gentry, S. Gorbunov, S. Halevi, V. Nikolaenko, G. Segev, V. Vaikuntanathan, and D. Vinayagamurthy, “Fully Key-Homomorphic Encryption, Arithmetic Circuit ABE, and Compact Garbled Circuits,” *IACR*, 2014.
- [8] D. Boneh, C. Gentry, S. Halevi, F. Wang, and D. J. Wu, “Private Database Queries using Somewhat Homomorphic Encryption,” in *ACNS*, 2013.
- [9] J. Fan and F. Vercauteren, “Somewhat Practical Fully Homomorphic Encryption,” *Cryptology ePrint Archive*, 2012.
- [10] C. Moore, M. O’Neill, E. O’Sullivan, Y. Doröz, and B. Sunar, “Practical Homomorphic Encryption: A Survey,” in *ISCAS 2014*.
- [11] P. Chaudhary, R. Gupta, A. Singh, and P. Majumder, “Analysis and Comparison of Various Fully Homomorphic Encryption Techniques,” in *GUCON*, 2013.
- [12] L. Organic, Y.-J. Chen, S. Dumas Ang, R. Lopez, X. Liu, K. Strauss, and L. Ceze, “Probing the Physical Limits of Reliable DNA Data Retrieval,” *Nature Communications*, 2020.
- [13] D. S. Cali, G. S. Kalsi, Z. Bingöl, C. Firtina, L. Subramanian, J. S. Kim, R. Ausavarungnirun, M. Alser, J. Gomez-Luna, A. Boroumand et al., “GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis,” in *MICRO*, 2020.
- [14] R. Bhukya and D. Somayajulu, “Exact Multiple Pattern Matching Algorithm using DNA Sequence and Pattern Pair,” *IJCA*, 2011.
- [15] C. Firtina, N. Mansouri Ghiasi, J. Lindegger, G. Singh, M. B. Cavlak, H. Mao, and O. Mutlu, “RawHash: Enabling Fast and Accurate Real-time Analysis of Raw Nanopore Signals for Large Genomes,” *Bioinformatics*, 2023.
- [16] E. De Cristofaro, S. Faber, and G. Tsudik, “Secure Genomic Testing with Size-and Position-Hiding Private Substring Matching,” in *WPES*, 2013.
- [17] M. M. A. Aziz, M. T. M. Tamal, and N. Mohammed, “Secure Genomic String Search with Parallel Homomorphic Encryption,” *Information*, 2024.
- [18] G. Verma and R. Chakraborty, “A Hybrid Privacy Preserving Scheme Using Finger Print Detection in Cloud Environment,” *ISI*, 2019.
- [19] W. Chen and Y. Gao, “Face Recognition using Ensemble String Matching,” *TIP*, 2013.
- [20] M. M. P. Mr, C. A. Dhote, and D. H. S. Mr, “Homomorphic Encryption for Security of Cloud Data,” *PROCS*, 2016.
- [21] A. Silberschatz, H. F. Korth, and S. Sudarshan, *Database System Concepts*, 1997.
- [22] N. Koudas, A. Marathe, and D. Srivastava, “Flexible String Matching against Large Databases in Practice,” in *VLDB*, 2004.
- [23] R. J. S. Raj, M. V. Prakash, T. Prince, K. Shankar, V. Varadarajan, and F. Nonyelu, “Web Based Database Security in Internet of Things Using Fully Homomorphic Encryption and Discrete Bee Colony Optimization,” *MJCS*, 2020.
- [24] A. A. Shaikh, “Attacks on Cloud Computing and its Countermeasures,” in *SCOPES*, 2016.
- [25] A. Duncan, S. Creese, and M. Goldsmith, “An Overview of Insider Attacks in Cloud Computing,” *CCPE*, 2015.
- [26] H. K. Bella and S. Vasundra, “A Study of Security Threats and Attacks in Cloud Computing,” in *ICSSIT*, 2022.
- [27] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshiba, “Secure Pattern Matching Using Somewhat Homomorphic Encryption,” in *CCSW*, 2013.
- [28] A. Essex, “Secure Approximate String Matching for Privacy-Preserving Record Linkage,” *TIFS*, 2019.
- [29] C. Bonte and I. Iliashenko, “Homomorphic String Search with Constant Multiplicative Depth,” in *CCS*, 2020.
- [30] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshiba, “Privacy-Preserving Wildcards Pattern Matching Using Symmetric Somewhat Homomorphic Encryption,” in *ACISP*, 2014.
- [31] A. Feer, “Privacy Preserving String Search using Homomorphic Encryption,” Master’s thesis, 2024.
- [32] A. Feldmann, N. Samardzic, A. Krastev, S. Devadas, R. Dreslinski, K. Eldefrawy, N. Genise, C. Peikert, and D. Sanchez, “F1: A Fast and Programmable Accelerator for Fully Homomorphic Encryption,” *MICRO 2021*.

- [33] G. Pradel and C. Mitchell, "Privacy-Preserving Biometric Matching Using Homomorphic Encryption," in *TrustCom*, 2021.
- [34] M. Kim, H. T. Lee, S. Ling, B. H. M. Tan, and H. Wang, "Private Compound Wildcard Queries using Fully Homomorphic Encryption," *TDSC*, 2017.
- [35] A. Bookstein, V. A. Kulyukin, and T. Raita, "Generalized Hamming Distance," *IR*, 2002.
- [36] S. Kim, J. Kim, M. J. Kim, W. Jung, J. Kim, M. Rhu, and J. H. Ahn, "BTS: An Accelerator for Bootstrappable Fully Homomorphic Encryption," in *ISCA 2022*.
- [37] N. Samardžić, A. Feldmann, A. Krastev, N. Manohar, N. Genise, S. Devadas, K. Eldefrawy, C. Peikert, and D. Sanchez, "Craterlake: A Hardware Accelerator for Efficient Unbounded Computation on Encrypted Data," in *ISCA 2022*.
- [38] D. E. Petrean and R. Potolea, "Homomorphic Encrypted Yara Rules Evaluation," *JISA*, 2024.
- [39] H. Gupta, M. Kabra, J. Gómez-Luna, K. Kanellopoulos, and O. Mutlu, "Evaluating Homomorphic Operations on a Real-World Processing-In-Memory System," in *IISWC*, 2023.
- [40] O. Mutlu, S. Ghose, J. Gómez-Luna, R. Ausavarungnirun, M. Sadrosadati, and G. F. Oliveira, "A Modern Primer on Processing-In-Memory," in *Emerging Computing: From Devices to Systems: Looking Beyond Moore and Von Neumann*, 2022.
- [41] J. Nider, C. Mustard, A. Zoltan, and A. Fedorova, "Processing in Storage Class Memory," in *HotStorage*, 2020.
- [42] A. Farmahini-Farahani, J. H. Ahn, K. Morrow, and N. S. Kim, "DRAMA: An Architecture for Accelerated Processing-near-Memory," *CAL*, 2014.
- [43] G. Singh, L. Chelini, S. Corda, A. J. Awan, S. Stuijk, R. Jordans, H. Corporaal, and A. J. Boonstra, "A Review of Near-Memory Computing Architectures: Opportunities and Challenges," in *Euromicro DSD*, 2018.
- [44] Y. Kwon, Y. Lee, and M. Rhu, "TensorDIMM: A Practical Near-memory Processing Architecture for Embeddings and Tensor Operations in Deep Learning," in *MICRO*, 2019.
- [45] S. Gupta, R. Cammarota, and T. Šimunić, "MemFHE: End-to-end Computing with Fully Homomorphic Encryption in Memory," *TECS*, 2024.
- [46] S. F. Yitbarek, T. Yang, R. Das, and T. Austin, "Exploring Specialized Near-memory Processing for Data-intensive Operations," in *DATE*, 2016.
- [47] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi, "A Scalable Processing-In-Memory Accelerator for Parallel Graph Processing," in *ISCA*, 2015.
- [48] J. Ahn, S. Yoo, O. Mutlu, and K. Choi, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-In-Memory Architecture," *ISCA*, 2015.
- [49] N. Hajinazar, G. F. Oliveira, S. Gregorio, J. D. Ferreira, N. M. Ghiasi, M. Patel, M. Alser, S. Ghose, J. Gómez-Luna, and O. Mutlu, "SIMDRAM: A Framework for Bit-Serial SIMD Processing Using DRAM," in *ASPLOS*, 2021.
- [50] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, "Ambit: In-memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," in *MICRO*, 2017.
- [51] J. D. Ferreira, G. Falcao, J. Gómez-Luna, M. Alser, L. Orosa, M. Sadrosadati, J. S. Kim, G. F. Oliveira, T. Shahroodi, A. Nori *et al.*, "pLUTO: Enabling Massively Parallel Computation in DRAM via Lookup Tables," in *MICRO*, 2022.
- [52] Y. Wang, L. Orosa, X. Peng, Y. Guo, S. Ghose, M. Patel, J. S. Kim, J. G. Luna, M. Sadrosadati, N. M. Ghiasi *et al.*, "FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching," in *MICRO*, 2020.
- [53] B. Gu, A. S. Yoon, D.-H. Bae, I. Jo, J. Lee, J. Yoon, J.-U. Kang, M. Kwon, C. Yoon, S. Cho *et al.*, "Biscuit: A Framework For Near-Data Processing of Big Data Workloads," *ISCA*, 2016.
- [54] I. Jo, D.-H. Bae, A. S. Yoon, J.-U. Kang, S. Cho, D. D. G. Lee, and J. Jeong, "YourSQL: A High-Performance Database System Leveraging In-Storage Computing," *PVLDB*, 2016.
- [55] N. M. Ghiasi, J. Park, H. Mustafa, J. Kim, A. Olgun, A. Gollwitzer, D. S. Cali, C. Firtina, H. Mao, N. A. Alserr *et al.*, "GenStore: A High-Performance and Energy-Efficient In-Storage Computing System for Genome Sequence Analysis," in *ASPLOS*, 2022.
- [56] J. Wang, D. Park, Y.-S. Kee, Y. Papakonstantinou, and S. Swanson, "SSD In-Storage Computing for List Intersection," in *DaMoN*, 2016.
- [57] D. Park, J. Wang, and Y.-S. Kee, "In-storage Computing for Hadoop MapReduce Framework: Challenges and Possibilities," *ToC*, 2016.
- [58] Z. Ruan, T. He, and J. Cong, "INSIDER: Designing In-Storage Computing System for Emerging High-Performance Drive," in *USENIX ATC*, 2019.
- [59] J. Wang, D. Park, Y. Papakonstantinou, and S. Swanson, "SSD In-storage Computing for Search Engines," *ToC*, 2016.
- [60] J. Park, R. Azizi, G. F. Oliveira, M. Sadrosadati, R. Nadig, D. Novo, J. Gómez-Luna, M. Kim, and O. Mutlu, "Flash-Cosmos: In-flash Bulk Bitwise Operations Using Inherent Computation Capability of NAND-Flash Memory," in *MICRO*, 2022.
- [61] J. Chen, C. Gao, Y. Lu, Y. Zhang, and J. Shu, "Ares-Flash: Efficient Parallel Integer Arithmetic Operations Using NAND Flash Memory," in *MICRO*, 2024.
- [62] C. Gao, X. Xin, Y. Lu, Y. Zhang, J. Yang, and J. Shu, "Parabit: Processing Parallel Bitwise Operations in NAND-Flash Memory-based SSDs," in *MICRO*, 2021.
- [63] O. Regev, "On Lattices, Learning with Errors, Random Linear Codes, and Cryptography," in *STOC*, 2005.
- [64] V. Lyubashevsky, C. Peikert, and O. Regev, "On Ideal Lattices and Learning with Errors Over Rings," in *EUROCRYPT*, 2010.
- [65] N. Smart and F. Vercauteren, "Fully Homomorphic SIMD Operations," *DCC*, 2014.
- [66] W. Castryck, I. Iliashenko, and F. Vercauteren, "Homomorphic SIMD Operations: Single Instruction Much More Data," in *EUROCRYPT*, 2018.
- [67] Q. Qin, N. Zhang, L. Liu, S. Yin, and S. Wei, "Addition Circuit Optimization Using Carry-Lookahead and SIMD for Homomorphic Encryption," in *EDSSC*, 2019.
- [68] M. Albrecht, M. Chase, H. Chen, J. Ding, S. Goldwasser, S. Gorbunov, S. Halevi, J. Hoffstein, K. Laine, K. Lauter, S. Lokam, D. Micciancio, D. Moody, T. Morrison, A. Sahai, and V. Vaikuntanathan, "Homomorphic Encryption Security Standard," Tech. Rep., 2018.
- [69] M. Olson, D. Davis, and J. W. Lee, "An Approach to the Exact Packed String Matching Problem," in *NLPIR*, 2020.
- [70] S. Sedghi, P. Van Liesdonk, S. Nikova, P. Hartel, and W. Jonker, "Searching Keywords with Wildcards on Encrypted Data," in *SCN*, 2010.
- [71] P. Kim, E. Jo, and Y. Lee, "An Efficient Search Algorithm for Large Encrypted Data by Homomorphic Encryption," *Electronics*, 2021.
- [72] M. Alser, Z. Bingöl, D. S. Cali, J. Kim, S. Ghose, C. Alkan, and O. Mutlu, "Accelerating Genome Analysis: A Primer on An Ongoing Journey," *IEEE Micro*, 2020.
- [73] M. Alser, J. Lindegger, C. Firtina, N. Almadhoun, H. Mao, G. Singh, J. Gomez-Luna, and O. Mutlu, "From Molecules to Genomic Variations: Accelerating Genome Analysis via Intelligent Algorithms and Architectures," *CSBJ*, 2022.
- [74] A. C. Yao, "Protocols for Secure Computations," in *SFCS*, 1982.
- [75] P. Snyder, "Yao's Garbled Circuits: Recent Directions and Implementations," in *Literature Review*, 2014.
- [76] O. Goldreich, S. Micali, and A. Wigderson, "How to Play any Mental Game, or A Completeness Theorem for Protocols with Honest Majority," in *Providing Sound Foundations for Cryptography: On the Work of Shafi Goldwasser and Silvio Micali*, 2019.
- [77] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, "Private Information Retrieval," *JACM*, 1998.
- [78] O. Goldreich, "Towards a Theory of Software Protection and Simulation by Oblivious RAMs," in *STOC*, 1987.
- [79] H. Chen, Z. Huang, K. Laine, and P. Rindal, "Labeled PSI from Fully Homomorphic Encryption with Malicious Security," in *CCS*, 2018.
- [80] M. Chase and E. Shen, "Substring-Searchable Symmetric Encryption," *PoPET*, 2015.
- [81] I. Leontiadis and M. Li, "Storage Efficient Substring Searchable Symmetric Encryption," in *CCS*, 2018.
- [82] G. S. Poh, J.-J. Chin, W.-C. Yau, K.-K. R. Choo, and M. S. Mohamad, "Searchable Symmetric Encryption: Designs and Challenges," *CSUR*, 2017.
- [83] M. Giraud, A. Anzala-Yamajako, O. Bernard, and P. Lafourcade, "Practical Passive Leakage-Abuse Attacks against Symmetric Searchable Encryption," in *SECURITY*, 2017.
- [84] P. Grubbs, K. Sekniqi, V. Bindschadler, M. Naveed, and T. Ristenpart, "Leakage-Abuse Attacks against Order-Revealing Encryption," in *SP*, 2017.
- [85] I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène, "TFHE: Fast Fully Homomorphic Encryption over the Torus," *JoC*, 2020.
- [86] Zama, "TFHE-rs: A Pure Rust Implementation of the TFHE Scheme for Boolean and Integer Arithmetics Over Encrypted Data," 2022, <https://github.com/zama-ai/tfhe-rs>.
- [87] M. Yasuda, T. Shimoyama, J. Kogure, K. Yokoyama, and T. Koshiha, "New Packing Method in Somewhat Homomorphic Encryption and Its Applications," *SCN*, 2015.
- [88] R. Micheloni, L. Crippa, and A. Marelli, *Inside NAND Flash Memories*, 2010.
- [89] R. Micheloni, A. Marelli, and K. Eshghi, *Inside Solid State Drives (SSDs)*, 2018.
- [90] T. Cho, Y.-T. Lee, E.-C. Kim, J.-W. Lee, S. Choi, S. Lee, D.-H. Kim, W.-G. Han, Y.-H. Lim, J.-D. Lee *et al.*, "A Dual-Mode NAND Flash Memory: 1-GB Multilevel and High-Performance 512-MB Single-Level Modes," *JSSC*, 2001.
- [91] S. Lee, Y.-T. Lee, W.-K. Han, D.-H. Kim, M.-S. Kim, S.-H. Moon, H. C. Cho, J.-W. Lee, D.-S. Byeon, Y.-H. Lim *et al.*, "A 3.3 V 4 Gb Four-Level NAND Flash Memory with 90 nm CMOS Technology," in *ISSCC*, 2004.
- [92] H. Maejima, K. Kanda, S. Fujimura, T. Takagiwa, S. Ozawa, J. Sato, Y. Shindo, M. Sato, N. Kanagawa, J. Musha *et al.*, "A 512Gb 3B/Cell 3D Flash Memory on a 96-Word-Line-Layer Technology," in *ISSCC*, 2018.
- [93] W. Cho, J. Jung, J. Kim, J. Ham, S. Lee, Y. Noh, D. Kim, W. Lee, K. Cho, K. Kim *et al.*, "A 1-TB, 4B/cell, 176-Stacked-WL 3D-NAND Flash Memory with Improved Read Latency and a 14.8 Gb/mm² Density," in *ISSCC*, 2022.
- [94] Y. Cai, G. Yalcin, O. Mutlu, E. F. Haratsch, A. Cristal, O. S. Unsal, and K. Mai, "Flash Correct-and-Refresh: Retention-aware Error Management for Increased Flash Memory Lifetime," in *ICCD*, 2012.
- [95] Y. Cai, Y. Luo, S. Ghose, and O. Mutlu, "Read Disturb Errors in MLC NAND Flash Memory: Characterization, Mitigation, and Recovery," in *DSN*, 2015.
- [96] Y. Cai, Y. Luo, E. F. Haratsch, K. Mai, and O. Mutlu, "Data Retention in MLC NAND Flash Memory: Characterization, Optimization, and Recovery," in *HPCA*, 2015.
- [97] Y. Cai, O. Mutlu, E. F. Haratsch, and K. Mai, "Program Interference in MLC NAND Flash Memory: Characterization, Modeling, and Mitigation," in *ICCD*, 2013.
- [98] Y. Cai, E. F. Haratsch, O. Mutlu, and K. Mai, "Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis, and Modeling," in *DATE*, 2013.
- [99] Q. Li, L. Shi, C. Gao, Y. Di, and C. J. Xue, "Access Characteristic Guided Read and Write Regulation on Flash-Based Storage Systems," *TOC*, 2018.
- [100] Y. Cai, S. Ghose, Y. Luo, K. Mai, O. Mutlu, and E. F. Haratsch, "Vulnerabilities in MLC NAND Flash Memory Programming: Experimental Analysis, Exploits, and Mitigation Techniques," in *HPCA*, 2017.
- [101] A. Tavakkol, M. Sadrosadati, S. Ghose, J. Kim, Y. Luo, Y. Wang, N. M. Ghiasi, L. Orosa, J. Gómez-Luna, and O. Mutlu, "FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives," in *ISCA*, 2018.
- [102] Samsung, "Samsung SSD 980 PRO," 2020. [Online]. Available: <https://www.samsung.com/semiconductor/minisite/ssd/product/consumer/980pro/>
- [103] S. A. R. 3.1, "Serial ATA International Organization," https://sata-io.org/system/files/specifications/SerialATA_Revision_3_1_Gold.pdf, 2011.
- [104] I. NVMe Express, "NVMe Express Base Specification, Revision 2.0d," <https://nvmexpress.org/specifications/>, 2021.
- [105] O. Mutlu, S. Ghose, J. Gómez-Luna, and R. Ausavarungnirun, "Processing Data Where It Makes Sense: Enabling In-Memory Computation," *MICPRO*, 2019.
- [106] N. M. Ghiasi, M. Sadrosadati, H. Mustafa, A. Gollwitzer, C. Firtina, J. Eudine, H. Mao, J. Lindegger, M. B. Cavlak, M. Alser, J. Park, and O. Mutlu, "MegIS:

- High-Performance, Energy-Efficient, and Low-Cost Metagenomic Analysis with In-Storage Processing,” in *ISCA*, 2024.
- [107] B. Yang, W. Xue, T. Zhang, S. Liu, X. Ma, X. Wang, and W. Liu, “End-to-End I/O Monitoring on Leading Supercomputers,” *TOS*, 2023.
- [108] R. Wong, N. Kim, K. Higgs, S. Agarwal, E. Ipek, S. Ghose, and B. Feinberg, “TCAM-SSD: A Framework for Search-Based Computing in Solid-State Drives,” *arXiv*, 2024.
- [109] D. Fakhry, M. Abdelsalam, M. W. El-Kharashi, and M. Safar, “A Review on Computational Storage Devices and Near Memory Computing for High-Performance Applications,” *MMDCS*, 2023.
- [110] K. Park, Y.-S. Kee, J. M. Patel, J. Do, C. Park, and D. J. Dewitt, “Query Processing on Smart SSDs,” *DEB*, 2014.
- [111] Y. Kang, Y.-s. Kee, E. L. Miller, and C. Park, “Enabling Cost-Effective Data Processing with Smart SSD,” in *MSST*, 2013.
- [112] G. Koo, K. K. Matam, T. I. H. K. G. Narra, J. Li, H.-W. Tseng, S. Swanson, and M. Annavaram, “Summarizer: Trading Communication with Computing Near Storage,” in *MICRO*, 2017.
- [113] W. Shim and S. Yu, “GP3D: 3D NAND based In-memory Graph Processing Accelerator,” *JETCAS*, 2022.
- [114] C.-Y. Chan and Y. E. Ioannidis, “Bitmap Index Design and Evaluation,” in *SIGMOD*, 1998.
- [115] E. O’Neil, P. O’Neil, and K. Wu, “Bitmap Index Design Choices and their Performance Implications,” in *IDEAS*, 2007.
- [116] Y. Li and J. M. Patel, “WideTable: An Accelerator for Analytical Data Processing,” in *VLDB*, 2014.
- [117] Y. Li and J. M. Patel, “BitWeaving: Fast Scans for Main Memory Data Processing,” in *SIGMOD*, 2013.
- [118] B. Goodwin, M. Hopcroft, D. Luu, A. Clemmer, M. Curmei, S. Elnikety, and Y. He, “BitFunnel: Revisiting Signatures for Search,” in *SIGIR*, 2017.
- [119] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, and M. A. Kozuch, “RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization,” in *MICRO*, 2013.
- [120] V. Seshadri, K. Hsieh, A. Boroumand, D. Lee, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, “Fast Bulk Bitwise AND and OR in DRAM,” *CAL*, 2015.
- [121] Fastbit, “FastBit: An Efficient Compressed Bitmap Index Technology,” <https://sdm.lbl.gov/fastbit/>.
- [122] M.-C. Wu and A. P. Buchmann, “Encoded Bitmap Indexing for Data Warehouses,” in *ICDE*, 1998.
- [123] Z. Guz, M. Awasthi, V. Balakrishnan, M. Ghosh, A. Shayesteh, T. Suri, and S. Semiconductor, “Real-Time Analytics as the Killer Application for Processing-In-Memory,” *WoNDP*, 2014.
- [124] Redis, “Redis bitmaps,” <https://redis.io/docs/data-types/bitmaps/>.
- [125] M. Alser, H. Hassan, H. Xin, O. Ergin, O. Mutlu, and C. Alkan, “GateKeeper: A New Hardware Architecture for Accelerating Pre-alignment in DNA Short Read Mapping,” *Bioinformatics*, 2017.
- [126] J. Loving, Y. Hernandez, and G. Benson, “BitPAL: A Bit-Parallel, General Integer-Scoring Sequence Alignment Algorithm,” *Bioinformatics*, 2014.
- [127] H. Xin, J. Greth, J. Emmons, G. Pekhimenko, C. Kingsford, C. Alkan, and O. Mutlu, “Shifted Hamming Distance: A Fast and Accurate SIMD-Friendly Filter to Accelerate Alignment Verification in Read Mapping,” *Bioinformatics*, 2015.
- [128] J. S. Kim, D. Senol Cali, H. Xin, D. Lee, S. Ghose, M. Alser, H. Hassan, O. Ergin, C. Alkan, and O. Mutlu, “Catalina-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies,” *BMC Genomics*, 2018.
- [129] E. S. Lander, L. M. Linton, B. Birren, C. Nusbaum, M. C. Zody, J. Baldwin, K. Devon, K. Dewar, M. Doyle, W. Fitzhugh *et al.*, “Initial Sequencing and Analysis of the Human Genome,” *Nature*, 2001.
- [130] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, “Basic Local Alignment Search Tool,” *JMB*, 1990.
- [131] G. Myers, “A Fast Bit-Vector Algorithm for Approximate String Matching Based on Dynamic Programming,” *JACM*, 1999.
- [132] D. Reis, J. Takeshita, T. Jung, M. Niemier, and X. S. Hu, “Computing-in-Memory for Performance and Energy-Efficient Homomorphic Encryption,” *TVLSI*, 2020.
- [133] M. Zhou, Y. Nam, P. Gangwar, W. Xu, A. Dutta, C. Wilkerson, R. Cammarota, S. Gupta, and T. Rosing, “FHEmem: A Processing In-Memory Accelerator for Fully Homomorphic Encryption,” *TETC*, 2025.
- [134] M. Mwaiesela, J. Hari, P. Yuhala, J. Ménétrey, P. Felber, and V. Schiavoni, “Evaluating the Potential of In-Memory Processing to Accelerate Homomorphic Encryption: Practical Experience Report,” in *SRDS*, 2024.
- [135] Y. Yang, H. Lu, and X. Li, “Poseidon-NDP: Practical Fully Homomorphic Encryption Accelerator Based on Near Data Processing Architecture,” *TCAD*, 2023.
- [136] T. Suzuki and H. Yamana, “Designing In-Storage Computing for Low Latency and High Throughput Homomorphic Encrypted Execution,” in *ICBDA*, 2023.
- [137] C. Gentry, S. Halevi, and N. P. Smart, “Better Bootstrapping in Fully Homomorphic Encryption,” in *IACR*, 2012.
- [138] “Microsoft SEAL,” <https://www.microsoft.com/en-us/research/project/microsoft-seal/>.
- [139] J. Kim, S. Kim, J. Choi, J. Park, D. Kim, and J. H. Ahn, “SHARP: A Short-word Hierarchical Accelerator for Robust and Practical Fully Homomorphic Encryption,” in *ISCA*, 2023.
- [140] V. S. Mailthody, Z. Qureshi, W. Liang, Z. Feng, S. G. De Gonzalo, Y. Li, H. Franke, J. Xiong, J. Huang, and W.-m. Hwu, “DeepStore: In-Storage Acceleration for Intelligent Queries,” in *MICRO*, 2019.
- [141] Y. Cho, M. H. Kim, and J.-S. Lee, “Memory Device, Operation Method of Memory Device, and Page Buffer Included in Memory Device,” Patent, 2022.
- [142] C. Kim, D.-H. Kim, W. Jeong, H.-J. Kim, I. H. Park, H.-W. Park, J. Lee, J. Park, Y.-L. Ahn, J. Y. Lee, S.-B. Kim, H. Yoon, J. D. Yu, N. Choi, N. Kim, H. Jang, J. Park, S. Song, Y. Park, J. Bang, S. Hong, Y. Choi, M.-S. Kim, H. Kim, P. Kwak, J.-D. Ihm, D. S. Byeon, J.-Y. Lee, K.-T. Park, and K.-H. Kyung, “A 512-Gb 3-b/Cell 64-Stacked WL 3-D-NAND Flash Memory,” *JSSC*, 2018.
- [143] Y.-C. Chen, Y.-H. Chang, and T.-W. Kuo, “Search-In-Memory: Reliable, Versatile, and Efficient Data Matching in SSD’s NAND Flash Memory Chip for Data Indexing Acceleration,” *TCAD*, 2024.
- [144] C. Kim, J. Ryu, T. Lee, H. Kim, J. Lim, J. Jeong, S. Seo, H. Jeon, B. Kim, I. Lee, D. Lee, P. Kwak, S. Cho, Y. Yim, C. Cho, W. Jeong, K. Park, J.-M. Han, D. Song, K. Kyung, Y.-H. Lim, and Y.-H. Jun, “A 21 nm High Performance 64 Gb MLC NAND Flash Memory with 400 MB/s Asynchronous Toggle DDR Interface,” *JSSC*, 2012.
- [145] H. Cao, F. Liu, Q. Wang, Z. Du, L. Jin, and Z. Huo, “An Efficient Built-in Error Detection Methodology with Fast Page-oriented Data Comparison in 3D NAND Flash Memories,” *Electronics Letters*, 2022.
- [146] K. E. Batcher, “Bit-Serial Parallel Processing Systems,” *TOC*, 1982.
- [147] N. Ganapathy and C. Schimmel, “General Purpose Operating System Support for Multiple Page Sizes,” in *USENIX ATC*, 1998.
- [148] W. Cheong, C. Yoon, S. Woo, K. Han, D. Kim, C. Lee, Y. Choi, S. Kim, D. Kang, G. Yu *et al.*, “A Flash Memory Controller for 15μs Ultra-Low-Latency SSD using High-Speed 3D NAND Flash with 3μs Read Time,” in *ISSCC*, 2018.
- [149] Intel Corp., “6th Generation Intel Core Processor Family Datasheet,” <http://www.intel.com/content/www/us/en/processors/core/>.
- [150] Intel, *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Vol. 3*, 2016.
- [151] S. Desrochers, C. Paradis, and V. M. Weaver, “A Validation of DRAM RAPL Power Measurements,” in *MEMSYS*, 2016.
- [152] K. N. Khan, M. Hirki, T. Niemi, J. K. Nurminen, and Z. Ou, “RAPL in Action: Experiences in Using RAPL for Power Measurements,” *TOMPECS*, 2018.
- [153] ARM, “ARM Cortex-R5,” <https://developer.arm.com/Processors/Cortex-R5>.
- [154] “QEMU,” https://docs.zephyrproject.org/latest/boards/qemu/cortex_r5/doc/index.html.
- [155] Micron Technology Inc., “4Gb: x4, x8, x16 DDR4 SDRAM Data Sheet,” 2016.
- [156] S. Ghose, T. Li, N. Hajinazar, D. S. Cali, and O. Mutlu, “Demystifying Complex Workload-DRAM Interactions: An Experimental Study,” *POMACS*, 2019.
- [157] G. Foundries, “Global Foundries Website,” <https://gf.com/technology-platforms/fdx-fd-soi/>, 2020.
- [158] B. Reidys, P. Liu, and J. Huang, “RSSD: Defend Against Ransomware with Hardware-Isolated Network-Storage Codesign and Post-Attack Analysis,” in *ASPLOS*, 2022.
- [159] Samsung, “Samsung. 2021. Samsung NVMe SSD 983 DCT: Enhanced for Speed and Reliability. White Paper,” https://download.semiconductor.samsung.com/resources/brochure/Data_Center_SSD_983_DCT_Product_Brief.pdf, 2021.
- [160] S. Faber, S. Jarecki, H. Krawczyk, Q. Nguyen, M. Rosu, and M. Steiner, “Rich Queries on Encrypted Data: Beyond Exact Matches,” in *ESORICS*, 2015.
- [161] F. Hahn, N. Loza, and F. Kerschbaum, “Practical and Secure Substring Search,” in *ICMD*, 2018.
- [162] N. Mainardi, A. Barengi, and G. Pelosi, “Privacy Preserving Substring Search Protocol with Polylogarithmic Communication Cost,” in *ACSAC*, 2019.
- [163] R. Guo, J. Li, and S. Yu, “{GridSE}: Towards Practical Secure Geographic Search via Prefix Symmetric Searchable Encryption,” in *USENIX Security*, 2024.
- [164] T. Moataz and A. Shikfa, “Boolean Symmetric Searchable Encryption,” in *SIGSAC*, 2013.
- [165] J. Daemen and V. Rijmen, “AES Proposal: Rijndael,” *NIST*, 1999.
- [166] D. Liestyowati, “Public Key Cryptography,” in *JPCS*, 2020.
- [167] T. M. Chan, “Optimal Partition Trees,” in *SOCG*, 2010.
- [168] F. Kerschbaum, “Frequency-Hiding Order-Preserving Encryption,” in *CCS*, 2015.
- [169] Z. Chen, Z. Cao, Z. Shen, and L. Ju, “HMC-FHE: A Heterogeneous Near Data Processing Framework for Homomorphic Encryption,” *TCAD*, 2024.