

Single-Ferroelectric FET based Associative Memory for Data-Intensive Pattern Matching

Abstract—Content addressable memories (CAMs) embeds parallel associative search directly into the memory blocks, thus finding widespread utility in associative memory (AM) related applications. To accommodate increasing demands of data-intensive search tasks, various efforts have been devoted to enhancing CAM density. These endeavors include the use of non-volatile memory (NVM) devices with compact structures and capitalizing on the multi-level cell (MLC) characteristics of NVM devices. In this work, we present a novel single-FeFET based CAM design, complemented by a 2-step search scheme. This design achieves ultra-compact storage density and supports dual CAM operations: binary/ternary CAM search for Hamming distance computations and multi-bit CAM for exact associative searches. Both binary/ternary CAM and multi-bit CAM operations have been illustrated and validated, and the area per bit, search latency and energy metrics have been evaluated at array level. In genome sequencing applications using hyperdimensional computing paradigm, our single-FeFET based AM engine achieves 89.9x/71.9x speedup and 66.5x/30.7x energy efficiency improvement over GPU implementations.

Index Terms—content addressable memory, ferroelectric field effect transistor, associative memory, pattern matching

I. INTRODUCTION

The exponentially growth in data generation and processing, driven by emerging computing models and tasks in edge devices and data centers, necessitates efficient computing hardware platforms. However, current mainstream digital computers, based on Von Neumann architectures, suffer from significant energy and performance costs due to the memory wall issues caused by substantial data transfer between memory and processors. Various solutions have been proposed to address the bottleneck, and one effective approach is compute-in-memory (CiM), which embeds memory-centric computations directly into memory blocks. This mitigates data movement, improving the performance and energy computing [1]–[9].

One type of CiM primitives employs content addressable memories (CAMs), which seamlessly integrate parallel associative search operations into the memory. CAMs find wide applications in associative memory (AM) related tasks such as inference and learning in artificial intelligence models [10]–[13] as shown in Fig. 1(a). They generate exact search results, determining whether stored words exactly match the query. Beyond exact search, CAMs have evolved to support approximate search, enabling parallel similarity computations between stored words and input query based on various distance metrics, such as Hamming distance and Euclidean distance [14]–[17]. This approximate search capability makes CAM as an ideal hardware for accelerating hyperdimensional computing (HDC) based pattern matching applications, including genome sequencing, image classification, speech cognition, etc. [18], [19]. HDC models encode classes into orthogonal hyper-vectors within high-dimension space, and transform complex

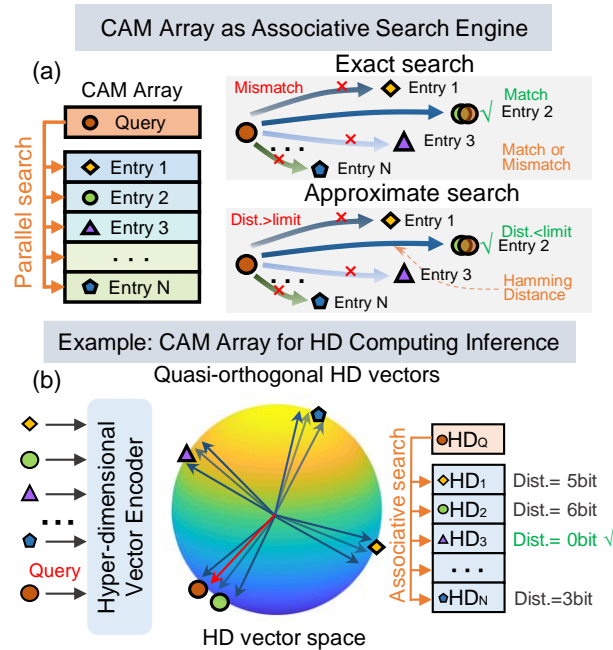


Fig. 1: Overview of CAM based associative search. (a) CAM array can perform exact search and approximate search in parallel. (b) An example of approximate search as associative memory in hyperdimensional computing.

sequential pattern matching into parallel Hamming distance computations. HDC inference involves computing Hamming distances between class hypervectors and input query (Fig. 1(b)), highly aligning with CAM engines [20], [21].

As data volumes continue to increase, high density CAM arrays are essential to accommodate data-intensive computing tasks. CAMs can be categorized based on the bit number of their stored value, including binary CAM (BCAM), ternary CAM (TCAM) with a wildcard state, multi-bit CAM (MCAM) for multi-bit values, and analog CAM (ACAM) for arbitrary value ranges [22]–[29]. Conventional CMOS based CAM designs typically utilize 10 transistors to store a binary value in static random access memory (SRAM) or 16 transistors to store a ternary value ('0', '1' and 'don't care') in an SRAM pair. More compact CAM designs have been built by leveraging the non-volatile storage properties of NVM devices, including two-terminal NVMs such as resistive random access memory (ReRAM), phase change memory (PCM), magnetic tunneling junctions (MTJ), and three-terminal devices like Spin-torque-transfer magnetic random access memory (STT-MRAM) and ferroelectric field effect transistor (FeFET). These designs frequently employ device pairs to store complementary bits or wildcard bits, significantly reducing the device count per cell. Recently, multi-level cell (MLC) characteristics of NVM devices have been explored to expand

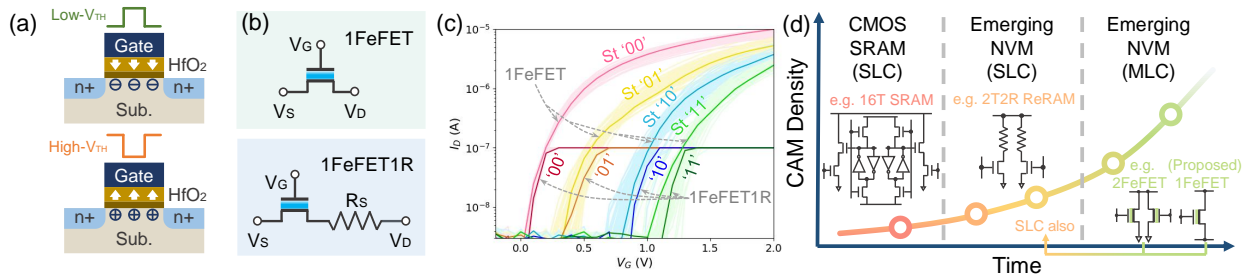


Fig. 2: (a) FeFET structure and polarization states after memory write pulses. (b) Schematic of 1FeFET and 1FeFET1R TCAM cells. (c) FeFET I_D - V_G curves for 60 different devices. The ON currents of 1FeFET1R cells are limited by resistors. (d) Existing CAM designs in terms of CAM bit density. The proposed 1FeFET based CAM design achieves the highest density.

beyond BCAM and TCAM designs. Notable examples include 6T-2R MCAM/ACAM utilizing the tunable conductance of ReRAMs [24] and FeFET based MCAM/ACAM utilizing the programmable threshold voltage states of FeFETs [25], [30].

That said, the aforementioned MCAM/ACAM designs typically require at least two active NVM devices, leaving the opportunity for denser and more energy efficient CAM designs unexplored. In this paper, we propose an ultra-compact CAM design that supports both BCAM and MCAM functions using just a single FeFET. We leverage the FeFET's single-transistor AND logic and voltage-driven MLC characteristics to achieve the dual functions. With our proposed 2-step search scheme, this single-FeFET based CAM design can function as a BCAM, performing parallel Hamming distance computations between stored entries and input query. Alternatively, it can operate as an MCAM, storing and searching multi-bit value for improved density. Such universal search-in-memory capability opens doors to applications in data-intensive tasks, particularly HDC based pattern matching for genome sequencing. We demonstrate the operation principles and functionality of this design, and also evaluate the area and performance metrics. Benchmarking results in HDC based genome sequencing tasks as AM engines show that our single-FeFET approach achieves 89.9x/71.9x speedup and 66.5x/30.7x higher energy efficiency compared to state-of-the-art DNA alignment tools.

II. BACKGROUND

In this section, we first introduce FeFET basics including the device and model, and then review existing BCAM and MCAM designs based on FeFET and other NVM devices.

A. FeFET Basics

FeFETs stand out as a promising emerging device candidate for embedded memory and low power CiM designs due to their CMOS-compatibility and efficient voltage driven read and write mechanisms. Fig. 2(a) illustrates the device structure that integrates HfO₂ ferroelectric (FE) layers into the gate stack. By applying positive or negative gate pulses to the device's gate, the ferroelectric polarization states within the FE layer are switched accordingly/partially. This switching behavior results in single-level cell (SLC) and MLC characteristics as depicted in Fig. 2(c). Moreover, recent innovations have included the integration of a series resistor at the drain of the FeFET (Fig. 2(b) [31]). This enhancement effectively mitigates device-to-device variability of FeFET ON current (Fig. 2(c)), substantially reducing sensing errors caused by

current variations. Therefore, we adopt 1FeFET-1R structure for our CAM design. The programmed MLC states can be read by applying a read bias and measuring the conduction current. Thanks to the voltage-driven write/read operation and three-terminal structure, which separates the write and read paths, FeFETs exhibit superior energy efficiency and more compact design methodology compared to other NVM devices [32].

We employ the Preisach FeFET model [33] which has been calibrated using experimental data for both device and circuit simulations. The Preisach model combines the responses of multiple FE domains and incorporates an underlying MOSFET model to characterize the hysteresis characteristics observed within partially polarized FE states. By applying gate voltage pulses with varying widths, FeFET model exhibits MLC curves calibrated by experiments, as shown in Fig. 2. In this paper, we leverage MLC FeFET model to build our proposed ultra-compact single FeFET CAM design.

B. Existing CAM Designs

Besides the matrix multiplications commonly used in neural networks, there's a fundamental need for parallel associative search and distance computations across memory blocks in various inference or leaning tasks. As a special type of CiM primitives, CAM designs have constantly evolved from binary CAM (BCAM) and ternary CAM (TCAM) to multi-bit CAM (MCAM) and analog CAM (ACAM), as depicted in Fig. 2(d). These CAM designs aim to enhance data storage density by leveraging the compact and MLC properties of non-volatile memory (NVM) devices. Most existing CAMs are BCAMs, which store binary values in a CAM cell and perform bit-wise XNOR logic to implement word-wise search operations. TCAMs, in addition to binary values, store a wildcard value, representing a 'don't care' state. Traditional CMOS based BCAM/TCAM cells comprise 10T/16T, incurring substantial power and area overheads [34], [35]. Compact and energy efficient BCAM and TCAM designs have been built using ReRAM and STT-MRAM devices. A 2T-2R TCAM design cell [36] utilizes binary states stored in memristor devices to represent logic values, while a 10T-4MTJ TCAM design [37] delivers rapid search operations. Nevertheless, both NVM based CAM designs still suffer from high write and search energy due to low variable resistance, low high-resistance-state (HRS)/low-resistance-state (LRS) ratios, and large access transistor or complex sensing circuitry needed for current-driven write and read mechanisms. FeFETs have emerged

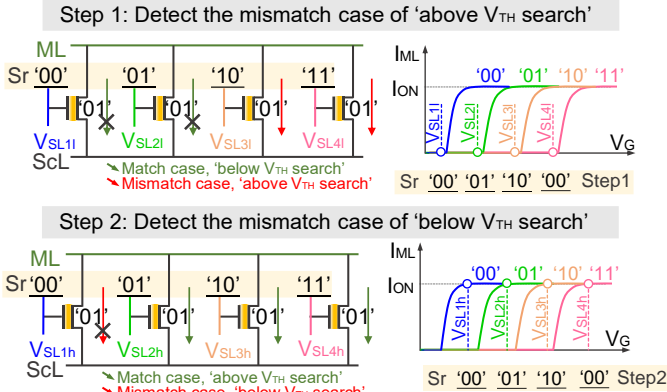


Fig. 4: 2-step search in the 1FeFET MCAM.

B. Multi-bit CAM for Exact Search

Our proposed single FeFET based CAM can also enable MCAM function, improving the data density. Similar to the search voltage and stored threshold voltage configurations, in the first step, the search voltages that are below V_{TH} identifies the mismatch cells with stored V_{TH} , and in the second step, the search voltages above V_{TH} identifies the mismatch cells with stored V_{TH} . Upon a match, ML currents in the two steps I_{MLS1}/I_{MLS2} will be low/high. Specifically, I_{MLS2} is roughly $I_{ON} \times N_{total}$. Such scheme enables only exact search functionality, while hamming distance computation is infeasible in MCAM. Nevertheless, the design is scalable depending on the number of distinct V_{TH} states FeFETs store.

Detailed operation principle of MCAM mode is illustrated as below, taking 2-bit per cell as an example. Without loss of generality, all cells store value '01' corresponding to the second lowest V_{TH} state, the search voltages encoding the query values '00', '01', '10', '11' are applied to the gates of FeFET cells, respectively, as shown in Fig. 4. In the first step, a search voltage i.e., V_{SL3l} , that is below the V_{TH} corresponding to value '10' and above the V_{TH} corresponding to value '01', is applied to search query value '10'. Then the mismatch cells that are applied by the search voltages above their stored V_{TH} states, i.e., the voltages searching for '10', '11', will conduct high ON currents I_{ON} . In the second step, for searching value '10', a search voltage, i.e., V_{SL3h} , above the V_{TH} corresponding to value '10', and below the V_{TH} corresponding to value '11', is applied. Then the mismatch cells that are applied by the search voltages below their stored V_{TH} states, i.e., the voltages searching for '00', will yield a low current, and all other cells will conduct high ON currents. The total ML current can determine whether such mismatch cells exist. As can be seen from above illustration, an exact match occurs only when the ML current is low in the first step, and linearly proportional to the number of cells, i.e., $I_{ML} \approx I_{ON} \times N_{total}$, in the second step. Other ML currents in the two steps indicate a mismatch condition. The proposed MCAM operation then completes.

C. Current-Domain Sense Amplifier

To support both BCAM and MCAM operations of our proposed single-FeFET based CAM, we adopt a current-domain thermometer-code analog-to-digital converter (ADC) as the

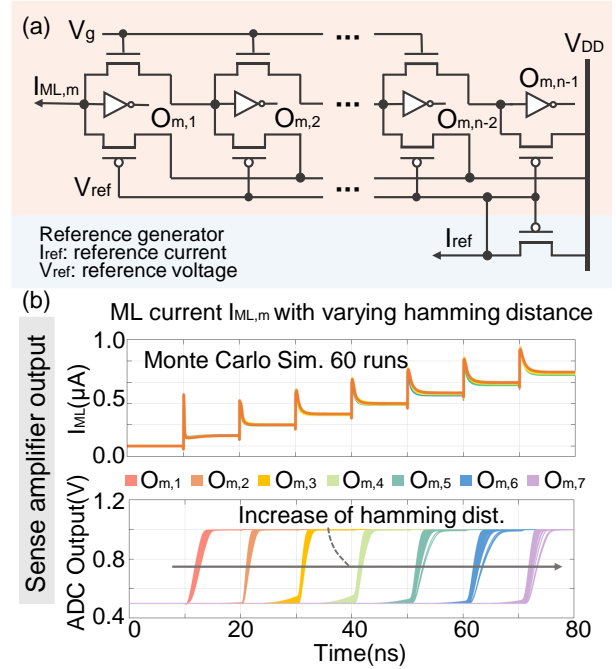


Fig. 5: (a) Schematic of thermometer-code ADC. (b) Output waveforms with varying Hamming distances.

sense amplifier (SA) to measure the ML currents in the 2-step scheme. The SA consists of a reference current generator for I_{ref} and a ladder-style current-to-voltage converter. Each converter of a word connects to its ML. The SA converts the ML current into voltage output as shown in Fig. 5. If the ML current I_{ML} from a CAM word is less than I_{ref} , the ladder-style current mirrors maintain high voltage at their drains, as the serial connected NMOS transistors conduct negligible current. If I_{ML} exceeds I_{ref} , the voltage at the first ladder current mirror's drain quickly drops down, activating the associated NMOS transistor and drawing excess current $I_{ML} - I_{ref}$. In this case, the output at the first ladder becomes high. As I_{ML} continues to grow and exceed two folds of I_{ref} , following the same principle, the output at the second ladder grows to high. Such ladder-style current-to-voltage conversion repeats for all ladders. Simulation results validate that the SA can detect the Hamming distance of the CAM word.

When the proposed design works in BCAM mode, the SA senses the currents of the 2-step scheme and calculate the Hamming distance between the input query and the stored words as discussed in Sec. III-A. When the design works in MCAM mode, the output of the SA directly indicates the match and mismatch results. We build the proposed CAM array with the SA for performance evaluations.

IV. EVALUATIONS

In this section, the proposed CAM design is evaluated and benchmarked in AM based HDC applications.

A. Function Validation

To validate the proposed single-FeFET based CAM design, simulations were conducted using a SPICE Monte Carlo method. Preisach FeFET model incorporating experimentally

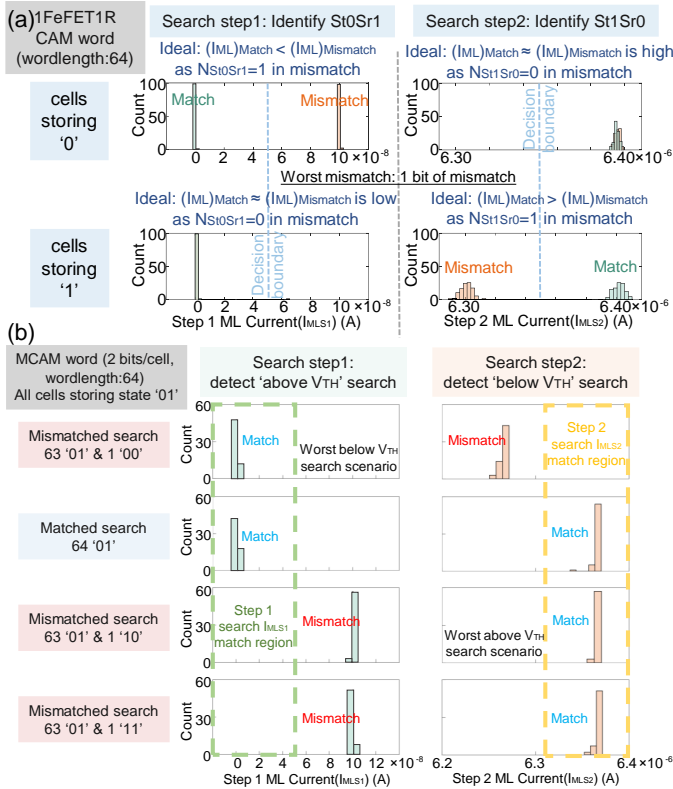


Fig. 6: 1FeFET CAM word functional verification. The ML current distributions of (a) a BCAM word and (b) a MCAM word indicate successful operations.

extracted device variations is employed for FeFETs, and predicted technology model (PTM) 45nm is used for MOSFETs [40]. An array with 64 wordlength is built for verification.

The Hamming distance computation of 1FeFET BCAM design is validated under two scenarios, where all cells store either '0' or '1'. The worst mismatch case is verified, where only one bit mismatch exists within the stored word. Per Fig. 6(a), the BCAM remains robust in the first search step. In the second step, thanks to the 1FeFET-1R structure, clear decision boundaries between different mismatch conditions can be defined to enable reliable Hamming distance computations.

We also validated the exact search function of the MCAM array with 64 cells per word. Without loss of generality, all cells store 2bit value '01', and the worst mismatch case is considered. Fig. 6(b) shows the output current distributions of the CAM word in the four cases. Only when all cells match with the input query, can the word ML currents are low in the first step and high in the second step. The mismatch cell upon a search value (i.e., '10'/'11') above the stored state (i.e., '01') results in an ON current in the first step, while the mismatch cell upon a search value (i.e., '00') below the stored state results in a less ML current in the second step. All these cases fail to meet the match conditions in the 2-step scheme.

B. Performance Evaluation

The single-FeFET based CAM design is evaluated in terms of area per bit, search latency and energy metrics. Fig. 7(a) shows the CAM density scaling trends based on various

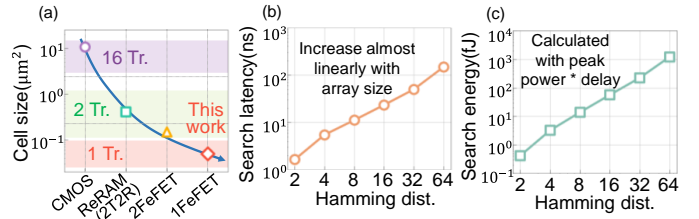


Fig. 7: (a) Cell size comparisons of 1FeFET CAM with CMOS, 2T-2R RAM and 2FeFET CAM designs. Search (b) latency and (c) energy of 1FeFET CAM array for one search step with different hamming distance thresholds.

emerging devices and CMOS technology. As can be seen, CAM schematic gradually evolves to 2T-2R and 2FeFET cells, which require only 2 devices. Our single-FeFET based CAM design ultimately provides the most compact cell and universal BCAM and MCAM functions. The MCAM mode further improves the density, achieving at least double density of the prior 2FeFET MCAM design [25].

Fig. 7(b)(c) demonstrates the search latency and energy of our CAM design regarding varying Hamming distances. The series resistor R_s is set to $1M\Omega$, and the search latency is defined as the time point when SA detects the maximum Hamming distance. The ladder-style ADC SA generates the voltage outputs in serial, thus the search latency and energy exhibit linearity to the number of ADC stages, i.e., the maximum Hamming distance or wordlength, which is consistent with Fig. 7. The latency can be reduced by adopting parallel ADC design, with more area and energy overheads [41].

C. Benchmarking on HDC

Our single-FeFET based CAM is further benchmarked in HDC genome sequencing tasks as AM kernels. Associative searches are widely used in genome sequencing, where a query DNA sequence comprising A, C, G, T nucleotide bases is searched across a long reference string to detect the presence of query sequence, and accelerate DNA alignment [42]. In the HDC architecture, the CAM array is integrated with TensorFlow and evaluated via a cycle-accurate simulator [43] given the array level performance. The overall HDC architecture is comprised of 32 tiles, each containing 128 512×512 CAM memory blocks. The genome database sequences from E.coli [44], Human CHR14 [44] and COVID-19 [45] are stored in the CAM blocks, and a genome query is searched across multiple CAMs in parallel. Hamming distance computations within a threshold are performed by CAMs to detect the closest reference sequences to the query. Fig. 8 suggests that our CAM based AM engine achieves on average 89.9X/71.9X and 66.5X/30.7X performance and energy efficiency improvements over state-of-the-art alignment tools NVBIO/GPU-BLAST [46]. This benchmarking highlights that our proposed single-FeFET based CAM design can be exploited for efficient AM engines, leveraging its high density and energy efficiency.

V. CONCLUSION

We propose a novel ultra-compact single-FeFET based CAM design that supports both BCAM and MCAM functions. The proposed 2-step search scheme encodes the binary and

