

A 3D MCAM architecture based on flash memory enabling binary neural network computing for edge AI

Maoying BAI¹, Shuhao WU¹, Hai WANG¹, Hua WANG¹, Yang FENG¹,
Yueran QI¹, Chengcheng WANG¹, Zheng CHAI², Tai MIN², Jixuan WU¹,
Xuepeng ZHAN^{1*} & Jiezhi CHEN^{1*}

¹*School of Information Science and Engineering, Shandong University, Qingdao 266200, China;*

²*Center for Spintronic and Quantum Systems, State Key Laboratory for Mechanical Behavior of Materials, School of Materials Science and Engineering, Xi'an Jiaotong University, Xi'an 710000, China*

Received 21 November 2023/Revised 22 February 2024/Accepted 24 April 2024/Published online 15 November 2024

Abstract The in-memory computing (IMC) architecture implemented by non-volatile memory units shows great possibilities to break the traditional von Neumann bottleneck. In this paper, a 3D IMC architecture is proposed whose unit is based on a multi-bit content-addressable memory (MCAM). The MCAM unit is comprised of two 65 nm flash memory and two transistors (2Flash2T), which is reconfigurable and multifunctional for both data write/search and XNOR logic operation. Moreover, the MCAM array can also support the population count (POPCOUNT) operation, which can be beneficial for the training and inference process in binary neural network (BNN) computing. Based on the well-known MNIST dataset, the proposed 3D MCAM architecture shows a 98.63% recognition accuracy and a 300% noise-tolerant performance without significant accuracy deterioration. Our findings can provide the potential for developing highly energy-efficient BNN computing for complex artificial intelligence (AI) tasks based on flash-based MCAM units.

Keywords reconfigurable, multifunction, multi-bit content-addressable memory (MCAM), bitwise operation, binary neural network, edge AI, flash memory, in-memory computing (IMC)

1 Introduction

In the era of big data, the increasing demand for efficient data processing has triggered an urgent challenge in the traditional von Neumann architecture, whose computing units and storage units are separated leading to a memory wall issue. With the capability to perform a storage function and logic function simultaneously, in-memory computing (IMC) is proposed and regarded as a potential candidate to break the von Neumann bottleneck [1,2]. There is a large amount of effort in developing IMC techniques both on hardware implementation and software algorithms [3–7]. At the hardware unit level, there is a lot of non-volatile memory including resistive random access memory (RRAM), ferroelectric random access memory (FeRAM), magnetoresistive random access memory (MRAM), phase change memory (PCM), and flash memory [8,9]. Owing to the merits of fast response to write/read operation and good adaptability to large arrays, flash memory has attracted extensive attention in the field of artificial intelligence at the edge (edge AI) including image recognition & classification, object detection & segmentation, and natural language processing [10–12]. Moreover, content-addressable memory (CAM) is one type of hardware unit that can easily search its entire contents in one clock cycle. With the unique feature of distinguishing the mismatching distance, CAM is capable of performing highly parallel and efficient search operations for data-intensive applications like pattern matching [13–16]. Compared to floating-point GPU implementations, Kazemi et al. [17] achieved similar accuracies by using flash-based multi bit-CAM (MCAM) with the ImageNet dataset, which significantly reduces energy consumption and operation latency. At the software algorithm level, inspired by the human nervous system, there are considerable numbers

* Corresponding author (email: zhanxuepeng@sdu.edu.cn, chen.jiezhi@sdu.edu.cn)

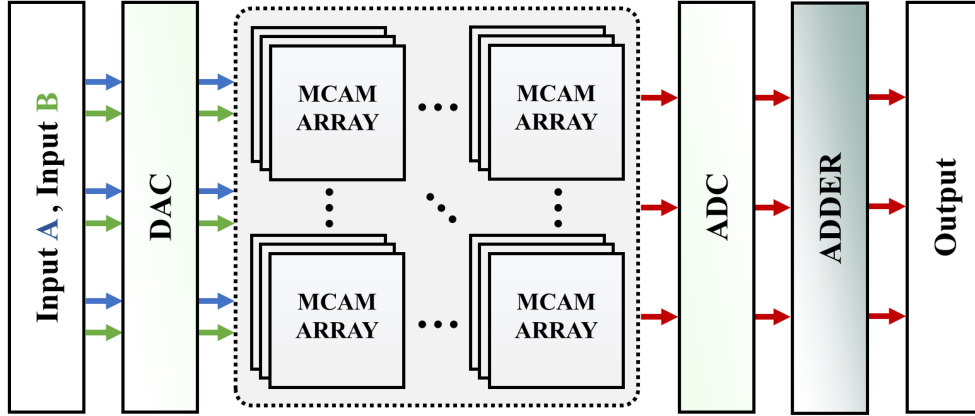


Figure 1 (Color online) Schematic image of the proposed 3D MCAM architecture with peripheral devices, which supports the binary-valued matrix multiplication.

of reports to construct highly-parallel and energy-efficient artificial neural networks (ANNs) including the deep neural network (DNN), convolutional neural network (CNN), binary neural network (BNN), recurrent neural network (RNN). Among them, BNN is a special one whose weights and activations are binary numbers (+1 or -1). Compared to other ANNs, BNN can support XNOR logic and population count (POPCOUNT) operations to supersede the floating-point operations for convolution and multiply-accumulate (MAC), which can significantly reduce the hardware resource [18–20]. Thus, BNN computing has received wide attention in various AI tasks in particular for those data-centered ones [21–24].

Although extensive research has been carried out on the employment of flashed-based CAM in BNN computing, there are still some challenges in reducing the computing hardware resource and realizing the offline training process. Firstly, at the unit level, several MCAM units and their peripheral circuits are required to realize the XNOR operation [25], which may increase memory space and energy costs. Secondly, there are normally two data lines (DL and \overline{DL}) at the traditional MCAM array for achieving MAC functions [26–28], which hinders the quick change/iterate the memory states and brings obstacles for off-line training. Thirdly, with the limited capability to distinguish the voltage in match lines (ML), the CAM arrays are preferred to certain AI tasks, which severely confines their further applications with complex datasets. Therefore, a compact, multifunctional, and reconfigurable flash-based MCAM unit is highly demanded, which may be of great importance for achieving highly energy-efficient and resource-constrained BNN computing with complex tasks.

In this paper, a reconfigurable MCAM unit (2Flash2T) supporting multi-bit stable states is proposed, which enables bitwise operations without additional peripheral devices. With a separated DL and \overline{DL} strategy, the proposed unit adopts a pair of input voltage vectors, which may be beneficial for achieving offline training in BNN computing. By mapping the partitioned matrixes to the multiple blocks, the proposed 3D MCAM architecture can perform typical AI tasks like image recognition, which shows an accuracy of $\sim 98.63\%$ with a noise-tolerant capability (input noise up to $\sim 300\%$) based on the Mixed National Institute of Standards and Technology database (MNIST). Our findings could provide a novel strategy to design 3D MCAM architecture for high energy-efficient, resource-constrained, and robust Edge AI based on the 65 nm flash memory.

2 MCAM unit and array characteristic

Figure 1 shows the schematic of the proposed 3D MCAM architecture for binary-valued matrix multiplications, which includes a digital-to-analog converter (DAC), 3D MCAM blocks, analog-to-digital converter (ADC), and adder subtractor (ADDER). The DAC is adapted to transfer the input matrixes to voltage matrixes, which are further transferred to binary matrixes by using ADC. The MCAM block (grey region) composed of several 3D MCAM arrays is used to realize the bitwise operations. The ADDER modules are utilized for integrating the output voltages from ADC leading to an output result. There are two binary-valued input matrixes (input \mathbf{A} and input \mathbf{B}) with the output corresponding to the result of matrix multiplication ($\mathbf{A} \times \mathbf{B}$). The bitwise matrix operations are mainly performed at the MCAM blocks with the help of peripheral modules.

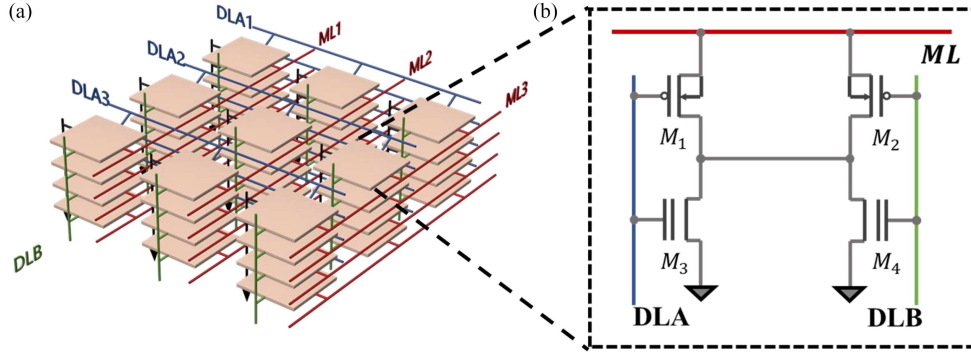


Figure 2 (Color online) Schematic image of (a) the 3D structure with the flash-based MCAM units and (b) the MCAM unit comprised of two depletion-type PMOS transistors (M_1 and M_2) and two flash memories (M_3 and M_4).

Table 1 Simulation setup key parameters and models

	Store -1	Store +1
M_1	$V_{TH-L} < V_{TH-P} < V_{TH-H}$	
M_2		
M_3	V_{TH-L}	V_{TH-H}
M_4	V_{TH-H}	V_{TH-L}
	Search -1	Search +1
DLA	V_{TH-L}	V_{TH-H}
DLB	V_{TH-H}	V_{TH-L}

Figure 2(a) shows the schematic of the 3D MCAM array, where the blue, green, and red lines stand for two data lines (DLA and DLB) and ML, respectively. The MCAM array contains $x \times y \times z$ MCAM units (Figure 2(b)), which are all comprised of two 65 nm flash memory and two depletion-type PMOS transistors (2Flash2T). Note that, x is the number of units paralleled in ML, and y is the column number of ML in the 3D array, and z is the layer number. Normally, two complementary DLs (DL and \overline{DL}) are used for single-datum operations in traditional CAM units. In our proposed MCAM unit, a separated DLs (DLA, DLB) strategy is adopted to represent 2 different data, which is achieved by the PMOS transistors to break the complementary constraint between two data lines. Note that in Figure 2, under the condition that apply a high voltage on the gate of PMOS and flash, the PMOS will not get broken. The reason is that the range of the programming and threshold voltage of flash memory is strictly controlled to fit the working voltage range of the PMOS. The line connecting the point between M_1 and M_3 to the point between M_2 and M_4 is necessary to ensure to CAM unit be able to work especially when searching “-1” of storing “+1” and searching “+1” of storing “-1”. Moreover, the setup parameters of the SPICE simulation are shown in Table 1 where $V_{TH-L} = 0.1$ V, $V_{TH-P} = 0.5$ V, $V_{TH-H} = 0.9$ V and the discharge capacitor connected between ML and the grounded is 0.1 pF.

Figure 3(a) displays the threshold voltage of individual flash memory and PMOS transistor simulated in SPICE. The 65 nm flash memory shows 16 storage states with threshold voltages (from V_{TH0} to V_{TH15}) ranging from 0.3 to 4.8 V. By adjusting the threshold voltages of flash memory to fit the threshold voltages of the PMOS (V_{TH-P}), the proposed unit can behave like typical MCAM and multifunctional MCAM. With V_{TH-P} larger than V_{TH15} (blue region), the proposed unit corresponds to a typical CAM unit with 16 different memory states. When V_{TH-P} is set within V_{TH0} and V_{TH15} (grey region), the proposed unit can also support the data write/search function with a reduced memory state. With V_{TH-P} smaller than V_{TH0} (green region), the proposed unit is always off (disabled). The 16 different transfer curves of the flash memory are shown in Figure 3(b) corresponding to state-0 ($V_{TH0} = 0.3$ V) to state-15 ($V_{TH15} = 4.8$ V). When mapping a matrix multiplication, the voltages on DLAs and DLBs represent the row vectors and column vectors in 2D scale respectively. In the 3D scale, the different DLAs and DLBs are mapped into different 2D planes. The benefits of the 3D structure lie in that the matrixes are only required to be mapped once to DLAs and DLBs. Owing to the separated DLs design, the mapping efficiency is improved, which is suitable for direct voltage iterating.

The PVT (process-voltage-temperature) analysis is conducted to reveal their impacts on the function of the CAM unit. The smaller discharge time (the time of V_{ML} reaches back to its saturated and stabilized

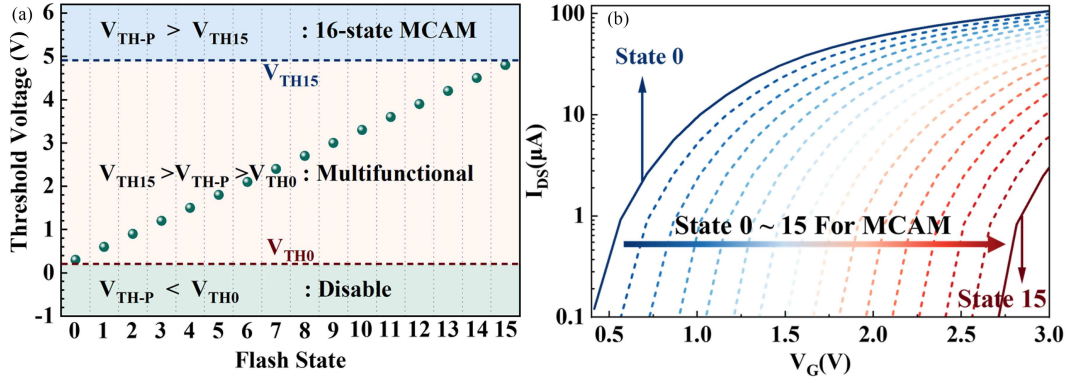


Figure 3 (Color online) (a) The threshold voltage distributions of a single flash and single PMOS; (b) the transfer curves of the flash memory with 16 different states corresponding to state-0 to state-15.

Table 2 Threshold voltage encoding of flash corresponding to different functions of the MCAM unit

$M_1 \& M_2$	$V_{TH-L} < V_{TH-P} < V_{TH-H}$			
Task	Case 1	Case 2	Case 3	Case 4
M_3	V_{TH-L}	V_{TH-L}	V_{TH-H}	V_{TH-H}
M_4	V_{TH-L}	V_{TH-H}	V_{TH-L}	V_{TH-H}
Truth table	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$
	$(-1, +1) = \text{Mismatch}$	$(-1, +1) = \text{Match}$	$(-1, +1) = \text{Mismatch}$	$(-1, +1) = \text{Match}$
	$(+1, -1) = \text{Mismatch}$	$(+1, -1) = \text{Mismatch}$	$(+1, -1) = \text{Match}$	$(+1, -1) = \text{Match}$
	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$	$(-1, -1) = \text{Match}$

value under matching condition) and larger voltage margin (the difference between the stable voltage of match and mismatch) can be obtained by using FF process corner, higher pre-charge voltage, and lower temperature [29–33], and larger voltage margin ratio can be obtained by lower pre-charge voltage. Compared to the process corner, the temperature and pre-charge voltage are dominant factors that affect the discharging time and voltage margin. Since the fabricating process is an important factor in affecting the flash V_t distribution [34–37], a V_t shift ranging from -0.2 to 0.2 V is adopted to reveal the impacts on the proposed CAM unit functionality in the SPICE simulation. The voltage margin decreases as the V_t shift increases. To guarantee a reliable function, the process variation induced V_t shift should be suppressed within 0.2 V for the proposed unit.

(a) The multifunctional MCAM Unit for write/search and XNOR logic functions. For the write/search function, the PMOS transistors always keep ON with $V_{TH-P} > V_{TH15}$, making the proposed MCAM unit behave like a typical MCAM unit with 16 memory states. Thus, the encoding scheme of the flash (M_3 and M_4) threshold voltage is the same as the traditional method.

For the XNOR logic operation, V_{TH-P} is required to be set to a certain range (the grey region in Figure 3(a)). The two binary data are transferred to input voltages imposed to DLA and DLB via DAC modules, which should be equal to the memory states (V_{TH0} to V_{TH15}). The lower input voltage is denoted as V_{TH-L} and the higher one is denoted as V_{TH-H} corresponding to logic -1 and logic $+1$, respectively. Note that the V_{TH-P} should be larger than V_{TH-L} and smaller than V_{TH-H} ($V_{TH-L} < V_{TH-P} < V_{TH-H}$).

As summarized in Table 2, different threshold voltages of M_3 and M_4 could lead to four different functions (Cases 1–4). Only when the two flashes are both programmed to V_{TH-L} , the XNOR function can be achieved. If V_{DLA} equals V_{DLB} , the MCAM unit could give rise to a match condition. Otherwise, a mismatch result is obtained, which corresponds to the XNOR logic function (Case 1). In the other cases (Cases 2–4), the truth tables are also displayed, which can be further adopted based on the user’s requirement. For example, Case 4 can be applied as wildcards to uniform the length when the data length is inconsistent. Moreover, with the separated DLs strategy, the input information of DLA and DLB (V_{DLA} and V_{DLB}) do not need to be saved in the MCAM unit, which is beneficial for rapid change.

(b) MCAM array characteristics for POPCOUNT functions. For the POPCOUNT logic operation, the MCAM array contains $x \times y \times z$ units that are adopted to count the number of mismatch conditions ((A, B) corresponds to $(-1, +1)$ and $(+1, -1)$). On the array level, x is limited by the saturated voltage on ML (V_{ML}) and its distinguish margin. Thus, in this work, 16 MCAM units are adopted to ensure

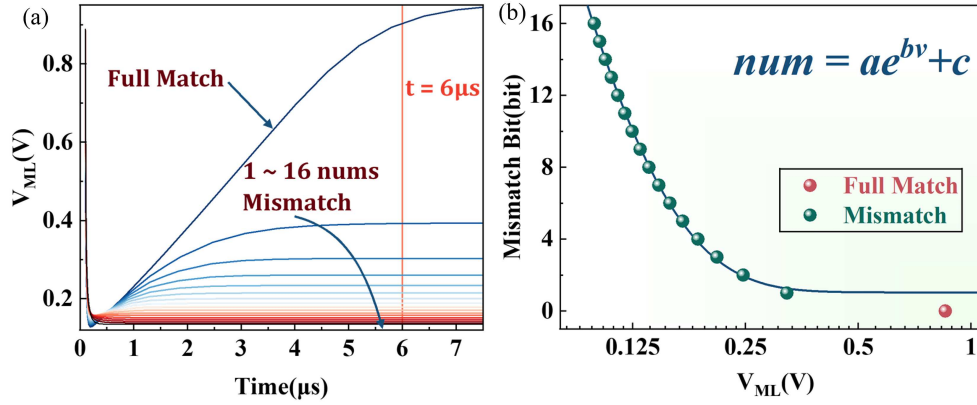


Figure 4 (Color online) (a) The discharge time-dependent distribution of ML current with various mismatching numbers; (b) the relationship of mismatching numbers versus the voltage of ML (V_{ML}) at the discharge time of 6 μ s.

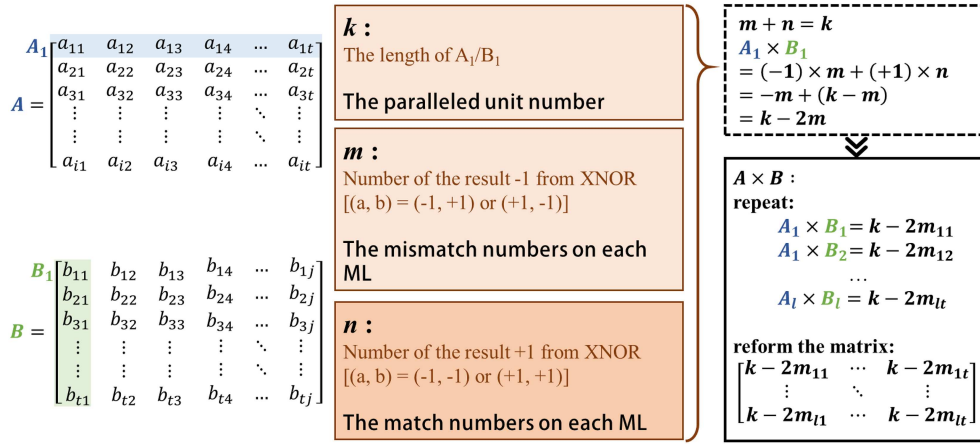


Figure 5 (Color online) Flow chart of the multiplication calculation process of $i \times t$ binary matrix A and $t \times j$ binary matrix B .

the mismatched conditions can be clearly distinguished. Figure 4(a) shows the relationship between discharge time and V_{ML} , which falls rapidly at first 0.5 μ s owing to the disturbance from the PMOS. If all 16 units are matched, the V_{ML} will recover to the recharging voltage, while the V_{ML} increases to various saturated voltages (smaller than the recharging voltage) corresponding to various mismatch numbers. With a discharge time of 6 μ s, the relationship between mismatching numbers versus the V_{ML} is displayed in Figure 4(b), which can be fitted as the following equation:

$$\text{num} = ae^{bv} + c, \quad (1)$$

where a , b , and c are all constant.

(c) MCAM architecture for binary matrix multiplication function. To perform binary matrix multiplication, periphery devices are further required in dealing with the POPCOUNT results.

Figure 5 shows the flow chart of the operating process for binary matrixes (A and B) multiplication, where A_1/B_1 stands for the first row/column vector. The number of XNOR results $-1/+1$ is assumed as m/n , which corresponds to the mismatch/match condition on ML. The length of the vector (A_1 and B_1) is denoted as k corresponding to the unit number on a single ML. The summary of mismatch conditions (m) and match conditions (n) equals (k):

$$m + n = k. \quad (2)$$

After the XNOR logic processing, the result of $A_1 \times B_1$ corresponds to

$$(+1) \times (k - m) + (-1) \times m, \quad (3)$$

which can be simplified to

$$k - 2m. \quad (4)$$

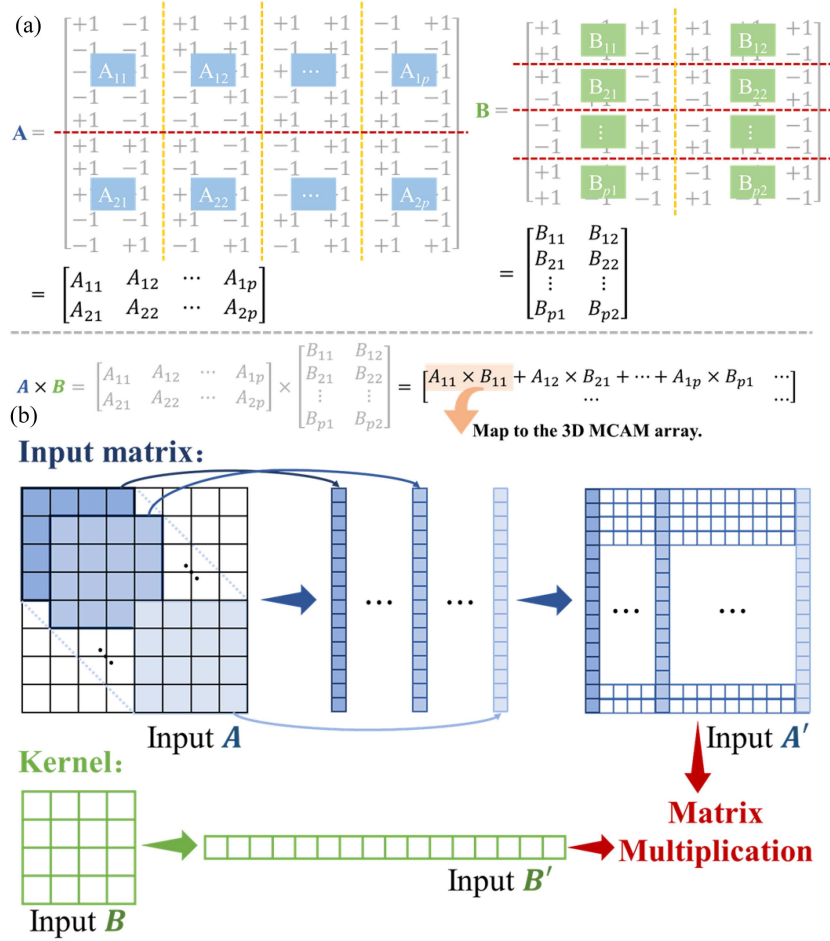


Figure 6 (Color online) Schematic image of (a) the matrix partitioning scheme for large-scale matrix multiplications and (b) converting convolution operations to matrix multiplications.

With the process repeated several times, the output of $A \times B$ can be obtained as illustrated in Figure 5. In this way, both the inputs for the matrix calculation can be applied to the circuit in the form of voltage, making rapid iteration for inputs possible.

3 3D MCAM array implement BNN computing

(a) The partitioned mapping for matrix multiplications and convolution operations. For typical AI tasks requiring large-scale matrix operations, a matrix partition method is adopted by using the proposed MCAM blocks. Figure 6(a) shows the schematic image for large-scale matrix multiplications. A and B are binary matrixes, which are segmented into smaller ones (A_{11} to A_{2p} and B_{11} to B_{p2}) that can be directly mapped to the MCAM arrays. The multiplication results of large matrixes can be obtained by successively integrating the partitioned matrixes.

Figure 6(b) shows the scheme of convolution operations with a 4×4 kernel (input B). The input matrix (input B) is unrolled step by step as the kernel slides upon it, leading to a recombinant matrix denoted as input A' . The convolution kernel is unrolled to a vector denoted as input B' . Then, the convolution results can be obtained by applying A' and B' into matrix multiplications.

(b) The network performances for BNN computing. To simplify the hardware implementation process, the BNN structure employed in this work is constructed by two convolutional layers and two fully connected layers. The first convolutional layer and the last fully-connected layer are performed with float point operations, while the others are binary. By using the well-known MNIST dataset, the network performances of BNN computing based on the proposed 3D MCAM architecture are evaluated. The recognition accuracy is mainly affected by two factors which are hyper-parameter setting and hardware

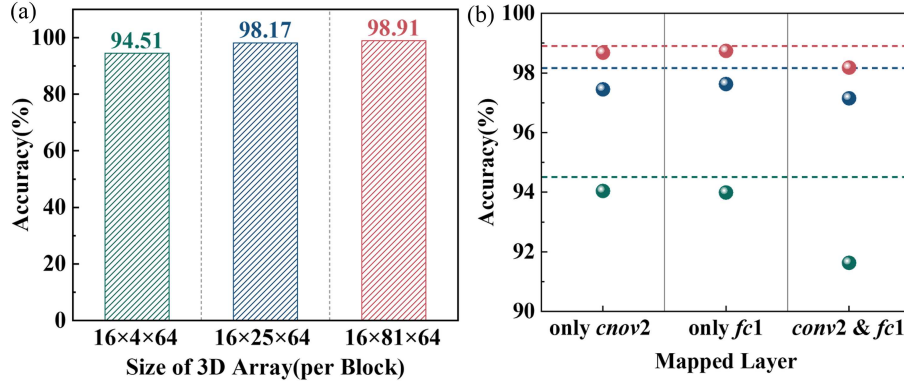


Figure 7 (Color online) (a) The different accuracies with various 3D array sizes in the form of $x \times y \times z$; (b) the accuracies versus mapped layer implemented by only cnov2 (first column), only fc1 (second column), as well as both cnov2 and fc1 (third column).

implementation. The hyper-parameter values correspond to different matrix scales in data flowing of BNN computing, which directly determines the array size. The training of the proposed CAM based BNN is achieved by using Bin_LeNet (the binarized LeNet) model. Several hyperparameters (the number of layers, convolutional kernel size, the i/o channel sizes) are adjusted to improve the network efficiency and accuracy. During training, there are some compressions like the convolution and pool operations for feature extraction, while no extra compression operations in the hardware simulation to maximize the accuracy.

Figure 7(a) shows the accuracies with different $x \times y \times z$ arrays (constant x & z of 16 & 64). Note that x is the number of units paralleled in ML, y is the column number of ML in the array, and z is the layer number of the output channel for cnov2. The accuracy increases slightly from 94.51% to 98.91% as y increases from 4 to 81. For different hardware implementation conditions, Figure 7(b) shows the accuracies with the mapped layer implemented by cnov2, only fc1, as well as both cnov2 and fc1 at various array sizes. It is clear that the hardware implementation has limited impacts on the accuracy. Further experimental results show that the circuit error has little effect on the training accuracy of large-scale neural networks. For more complex datasets, a large array size is normally required, which brings challenges in increasing the number of CAM units paralleled in one ML and of MLs paralleled in one plane. This may be achieved by suppressing the device variation and enlarging the voltage margin. Additionally, a dynamical matrix partition strategy is also important, which can adjust the matrix size based on different data complexities and required accuracies. The proposed architecture can be applied to some more complex datasets. For example, the classification accuracy is $\sim 78.09\%$ by using Bin_VGG13 model in the CIFAR10 dataset. The network performances can be optimized by using a complex neural network model (Bin_VGG16, Bin_VGG19) [37], enlarging the array size and improving the CAM unit performance.

Moreover, the noise immunity is verified based on the MCAM block implemented BNN computing. Figure 8 shows the accuracy (left axis) and error (right axis) under different noise rates (α), which is a coefficient of the imposed noise matrix. The noise matrix is comprised of random values (0–1) with the same size as the MNIST images which have pixel values ranging from 0 to 1 originally. It is clear that the accuracy decreases slightly until the noise rate α approaches 3.5, which indicates a strong noise-tolerant capability of the proposed MCAM-based BNN computing.

(c) Comparison with other CAM units and peripheral circuit. Table 3 [22,38–42] shows the comparison of different CAM units' performance with our work. The energy efficiency is defined as the average energy consumption per search for each unit, which is about 0.18 fJ/bit/search in our work. The latency is defined as the search delay between the rising edges of the clock and V_{ML} , which corresponds to the time interval of pre-charge stage and match stage [40]. The area is evaluated by the device structures & number and process node, considering the device dimensions are unknown in the SPICE simulation. Compared to other related work, the proposed CAM shows advantages in terms of average energy consumption, reconfigurable characteristics, multi-state storage and bitwise operation. The proposed architecture shows advantages in the multi-functional characteristics, high-integrated array and potential for off-line training. The architecture can accomplish the typical 16-state MCAM function and bitwise operation of XNOR calculations with energy consumption of 0.18 fJ/bit/search. Benefitting from the reliable characteristics

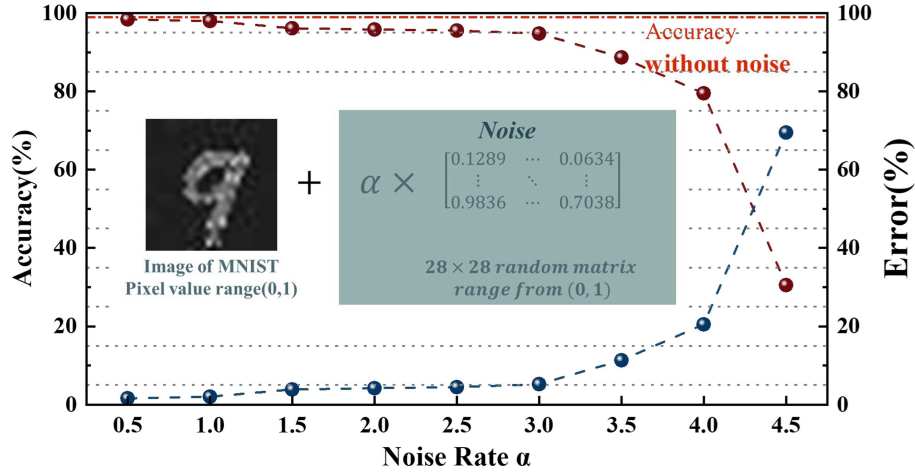


Figure 8 (Color online) Recognition accuracy (left) and error (right) with noise disturbance, where the noise is 28×28 random matrix between 0–1 and the rate α is a coefficient.

Table 3 Performance comparison by using different units

	This work	[22]	[38]	[39]	[40]	[41]	[42]
Area (TCAM cell)	2F2T	12T	2T2R	2.5T1R	2FeFET1T	2T2F	2F
Technology	65 nm+Flash	65 nm+SRAM	65 nm+RRAM	65 nm+RRAM	45 nm+FeFET	32 nm+FGFET	55 nm+Flash
Search delay (ns)	25	–	0.735	1	0.25	~20	2
Avg. E_{search} (fJ/bit/search)	0.18	2.48	0.17	0.495	0.195	–	0.34
Bitwise operation	Y	Y	–	–	–	Y	Y
Storage capacity	1–16	1	1	1	–	1	8
Reconfigurable	Y	N	N	N	N	N	N

and mature process of flash memory, the proposed CAM unit is easy for large-scale and 3D integration. The design of two different DLs may provide the possibility for accomplishing off-line training.

4 Conclusion

In this work, a reconfigurable MCAM unit (2Flash2T) with 14 stable states is proposed to realize the XNOR and POPCOUNT operations. This allows the input values applied on the reconfigurable unit to realize a 14-state search function as a traditional MCAM unit does. Based on the block matrix multiplication scheme, the proposed 3D MCAM array is capable of disposing of the well-known MNIST dataset and receives a recognition accuracy of up to $\sim 98.63\%$ and a noise-tolerant capability with $\sim 300\%$ input noise without apparent accuracy drop. A novel design of 3D MCAM architecture based on the 65 nm flash memory for high energy-efficient, resource-constrained, and robust edge AI tasks is provided by the findings of our work.

Acknowledgements This work was supported by National Key Research and Development Program of China (Grant Nos. 2023YFB4402500, 2023YFB4402400), National Natural Science Foundation of China (Grant Nos. 62034006, 92264201, U23B2040), Natural Science Foundation of Shandong Province (Grant Nos. ZR2023LZH007, TSQN202306059), and Program of Qilu Young Scholars of Shandong University and Micro-Nanoelectronic Device Project (Grant No. MINDXZ202407).

References

- Ielmini D, Wong H S P. In-memory computing with resistive switching devices. *Nat Electron*, 2018, 1: 333–343
- Li Y, Wang Z, Midya R, et al. Review of memristor devices in neuromorphic computing: materials sciences and device challenges. *J Phys D-Appl Phys*, 2018, 51: 503002
- Zidan M A, Strachan J P, Lu W D. The future of electronics based on memristive systems. *Nat Electron*, 2018, 1: 22–29
- Li C, Wang Z, Rao M, et al. Long short-term memory networks in memristor crossbar arrays. *Nat Mach Intell*, 2019, 1: 49–57
- Sebastian A, Le Gallo M, Khaddam-Aljameh R, et al. Memory devices and applications for in-memory computing. *Nat Nanotechnol*, 2020, 15: 529–544
- Feng Y, Chen B, Tang M, et al. Near-threshold-voltage operation in flash-based high-precision computing-in-memory to implement Poisson image editing. *Sci China Inf Sci*, 2023, 66: 222402

- 7 Zhan X, Chen J, Ji Z. Insights of VG-dependent threshold voltage fluctuations from dual-point random telegraph noise characterization in nanoscale transistors. *Sci China Inf Sci*, 2022, 65: 189405
- 8 Yu S, Jiang H, Huang S, et al. Compute-in-memory chips for deep learning: recent trends and prospects. *IEEE Circ Syst Mag*, 2021, 21: 31–56
- 9 Jeloka S, Akesh N B, Sylvester D, et al. A 28 nm configurable memory (TCAM/BCAM/SRAM) using push-rule 6T bit cell enabling logic-in-memory. *IEEE J Solid-State Circ*, 2016, 51: 1009–1021
- 10 Lue H T, Hsu P K, Wei M L, et al. Optimal design methods to transform 3D NAND flash into a high-density, high-bandwidth and low-power nonvolatile computing in memory (nvCIM) accelerator for deep-learning neural networks (DNN). In: *Proceedings of IEEE International Electron Devices Meeting (IEDM)*, 2019
- 11 Xiang Y, Huang P, Han R, et al. Efficient and robust spike-driven deep convolutional neural networks based on NOR flash computing array. *IEEE Trans Electron Dev*, 2020, 67: 2329–2335
- 12 Han R, Huang P, Xiang Y, et al. A novel convolution computing paradigm based on NOR flash array with high computing speed and energy efficiency. *IEEE Trans Circ Syst I*, 2019, 66: 1692–1703
- 13 Kwak B, Kim H, Kwon D. Ferroelectric-gate tunnel field-effect transistor one-transistor ternary contents addressable memory. *Semicond Sci Technol*, 2023, 38: 055013
- 14 Chen Y, Mu J, Kim H, et al. BP-SCIM: a reconfigurable 8T SRAM macro for bit-parallel searching and computing in-memory. *IEEE Trans Circ Syst I*, 2023, 70: 2016–2027
- 15 Jianwei Z, Yizheng Y, Binda L, et al. A cascaded charge-sharing technique for an EDP-efficient match-line design in CAMs. *J Semicond*, 2009, 30: 065009
- 16 Zhuo C, Yang Z, Ni K, et al. Design of ultracompact content addressable memory exploiting 1T-1MTJ cell. *IEEE Trans Comput-Aided Des Integr Circ Syst*, 2023, 42: 1450–1462
- 17 Kazemi A, Sahay S, Saxena A, et al. A flash-based multi-bit content-addressable memory with Euclidean squared distance. In: *Proceedings of IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*, 2021
- 18 Rastegari M, Ordonez V, Redmon J, et al. XNOR-Net: ImageNet classification using binary convolutional neural networks. 2016. ArXiv:1603.05279
- 19 Courbariaux M, Hubara I, Soudry D, et al. Binarized neural networks: training deep neural networks with weights and activations constrained to +1 or -1. 2016. ArXiv:1602.02830
- 20 Deng L, Jiao P, Pei J, et al. GXNOR-Net: training deep neural networks with ternary weights and activations without full-precision memory under a unified discretization framework. *Neural Netw*, 2018, 100: 49–58
- 21 Si X, Chang M F, Khwa W S, et al. A dual-split 6T SRAM-based computing-in-memory unit-macro with fully parallel product-sum operation for binarized DNN edge processors. *IEEE Trans Circ Syst I*, 2019, 66: 4172–4185
- 22 Yin S, Jiang Z, Seo J S, et al. XNOR-SRAM: in-memory computing SRAM macro for binary/ternary deep neural networks. *IEEE J Solid-State Circ*, 2020, 55: 1733–1743
- 23 Qiu H, Ma H, Zhang Z, et al. RBNN: memory-efficient reconfigurable deep binary neural network with IP protection for Internet of Things. *IEEE Trans Comput-Aided Des Integr Circ Syst*, 2023, 42: 1185–1198
- 24 Halawani Y, Mohammad B, Abu Lebdeh M, et al. ReRAM-based in-memory computing for search engine and neural network applications. *IEEE J Emerg Sel Top Circ Syst*, 2019, 9: 388–397
- 25 Chen Y, Lu L, Kim B, et al. Reconfigurable 2T2R ReRAM with split word-lines for TCAM operation and in-memory computing. In: *Proceedings of IEEE International Symposium on Circuits and Systems (ISCAS)*, 2020
- 26 Laguna A F, Yin X, Reis D, et al. Ferroelectric FET based in-memory computing for few-shot learning. In: *Proceedings of the Great Lakes Symposium on VLSI*, 2019. 373–378
- 27 Wang X, Wang L, Wang Y, et al. A 4T2R RRAM bit cell for highly parallel ternary content addressable memory. *IEEE Trans Electron Dev*, 2021, 68: 4933–4937
- 28 Sagario C J, Iii B Q, Jimenez K G, et al. Design of single poly flash memory cell with power reduction technique at program mode in 65nm CMOS process. In: *Proceedings of International Conference on Control, Electronics, Renewable Energy and Communications (ICCEREC)*, 2018
- 29 Resnati D, Goda A, Nicosia G, et al. Temperature effects in NAND flash memories: a comparison between 2-D and 3-D arrays. *IEEE Electron Device Lett*, 2017, 38: 461–464
- 30 Lee W, Park C, Kim K. Temperature dependence of endurance characteristics in NOR flash memory cells. In: *Proceedings of IEEE International Reliability Physics Symposium Proceedings*, 2006. 701–702
- 31 Wang J, Bai Y, Wang H, et al. Reconfigurable bit-serial operation using toggle SOT-MRAM for high-performance computing in memory architecture. *IEEE Trans Circ Syst I*, 2022, 69: 4535–4545
- 32 Ali M F, Jaiswal A, Roy K. In-memory low-cost bit-serial addition using commodity DRAM technology. *IEEE Trans Circ Syst I*, 2020, 67: 155–165
- 33 An H, Kim K, Jung S, et al. The threshold voltage fluctuation of one memory cell for the scaling-down NOR flash. In: *Proceedings of the 2nd IEEE International Conference on Network Infrastructure and Digital Content*, 2010. 433–436
- 34 Li H. Modeling of threshold voltage distribution in NAND flash memory: a Monte Carlo method. *IEEE Trans Electron Dev*, 2016, 63: 3527–3532
- 35 Kim W, Kim Y, Park S H, et al. Variation of threshold voltage and ON-cell current caused by cell gate length fluctuation in virtual source/drain NAND flash memory. *Jpn J Appl Phys*, 2012, 51: 074301
- 36 Yang T, Xia Z, Shi D, et al. Analysis and optimization of threshold voltage variability by polysilicon grain size simulation in 3D NAND flash memory. *IEEE J Electron Dev Soc*, 2020, 8: 140–144
- 37 Nugraha G S, Darmawan M I, Dwiyanaputra R. Comparison of CNN's architecture GoogleNet, AlexNet, VGG-16, Lenet-5, Resnet-50 in Arabic handwriting pattern recognition. *KINETIK*, 2023, 8: 2
- 38 Pan K, Tosson A M S, Wang N, et al. A novel cascaded TCAM using RRAM and current race scheme for high-speed energy-efficient applications. *IEEE Trans Nanotechnol*, 2023, 22: 214–221
- 39 Lin C C, Hung J Y, Lin W Z, et al. 7.4 A 256b-wordlength ReRAM-based TCAM with 1ns search-time and 14× improvement in wordlength-energyefficiency-density product using 2.5T1R cell. In: *Proceedings of IEEE International Solid-State Circuits Conference (ISSCC)*, 2016. 136–137
- 40 Yin X, Qian Y, Imani M, et al. Ferroelectric ternary content addressable memories for energy-efficient associative search. *IEEE Trans Comput-Aided Des Integr Circ Syst*, 2023, 42: 1099–1112
- 41 Cho S, Kim S, Choi I, et al. Non-volatile logic-in-memory ternary content addressable memory circuit with floating gate field effect transistor. *AIP Adv*, 2023, 13: 045211
- 42 Kazemi A, Sahay S, Saxena A, et al. A flash-based multi-bit content-addressable memory with Euclidean squared distance. In: *Proceedings of IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*, 2021