

Cite this: *Nanoscale Adv.*, 2023, 5, 1559

Research progress in architecture and application of RRAM with computing-in-memory

Chenyu Wang,^a Ge Shi,^{id}^a Fei Qiao,^{id}^b Rubin Lin,^a Shien Wu^a and Zenan Hu^a

The development of new technologies has led to an explosion of data, while the computation ability of traditional computers is approaching its upper limit. The dominant system architecture is the von Neumann architecture, with the processing and storage units working independently. The data migrate between them *via* buses, reducing computing speed and increasing energy loss. Research is underway to increase computing power, such as developing new chips and adopting new system architectures. Computing-in-memory (CIM) technology allows data to be computed directly on the memory, changing the current computation-centric architecture and designing a new storage-centric architecture. Resistive random access memory (RRAM) is one of the advanced memories which has appeared in recent years. RRAM can change its resistance with electrical signals at both ends, and the state will be preserved after power-down. It has potential in logic computing, neural networks, brain-like computing, and fused technology of sense-storage-computing. These advanced technologies promise to break the performance bottleneck of traditional architectures and dramatically increase computing power. This paper introduces the basic concepts of computing-in-memory technology and the principle and applications of RRAM and finally gives a conclusion about these new technologies.

Received 11th January 2023

Accepted 4th February 2023

DOI: 10.1039/d3na00025g

rsc.li/nanoscale-advances

1 Introduction

Scientific research and social life bring vast amounts of data. Researchers are constantly upgrading computing power to cope with large amounts of information processing tasks. Within the constraints of the von Neumann architecture, the dominant approach to increase the overall system speed is to increase the performance of chips. As semiconductor processing is approaching physical limits, Moore's Law is about to expire.^{1,2} The number of transistors per unit area of a chip will no longer increase significantly. The storage wall and power wall created by data migration have also limited the performance of computers. There has not been an order-of-magnitude increase in computing power.

Computing-in-memory technology integrates computing and storage capabilities into a single unit. Data need not migrate between the processor and memory. The technology breaks the bottleneck of traditional computer architecture, which is considered a significant development trend for future breakthroughs in computing power.³ W. H. Kautz⁴ of the Stanford Research Institute first proposed the concept of computing-in-memory in the 1970s but was unable to implement it due to technical constraints. Early computing-in-

memory solutions were usually based on mature CMOS processes. The research is focused on the existing memory cell circuits, including static random access memory (SRAM), dynamic random access memory (DRAM), and Flash.⁵ The technology can make logic operations, and the performance and cost of the device are reduced compared with those of the conventional component. However, with the advent of new non-volatile memories and the progress of semiconductor processing, computing-in-memory technology has been further developed.

Advanced storage technologies include resistive storage, ferroelectric storage,⁶ phase change storage,^{7,8} and magnetic storage.⁹ Resistive storage is a current research hotspot. The simple structure, high integration, low power consumption, and high speed make it one of the advantageous candidates for next-generation non-volatile memory. In 1971, the concept of memristors was first proposed by Professor Chua¹⁰ at the University of California, Berkeley. It was not until 2008 that a prototype of the memristor was successfully prepared by Hewlett-Packard Laboratories.¹¹ In the time that follows, research studies focus on mathematical models,¹² resistance mechanisms,^{13,14} across-array structures, preparation processing,¹⁵ and fabrication materials.^{16,17} The memristor used in memory applications is called resistive random-access memory (RRAM). It is considered the key to breaking through AI arithmetic and big data, to meet high-speed computing and also to meet low-precision, high-speed computing needs raised by IoT and neural network computing. The computing-in-memory

^aCollege of Mechanical and Electrical Engineering, China Jiliang University, Hangzhou, China. E-mail: shige@cjl.u.edu.cn

^bDept of Electronic Engineering, Tsinghua University, Beijing 310018, People's Republic of China



structure is expected to replace von Neumann architecture chips in many fields in the future.¹⁸

2 Computing-in-memory technology

The von Neumann architecture is characterized by the fact that the processing and storage units are independent, and the two are connected *via* a bus. However, traditional architecture has encountered some problems. On the one hand, data migration across the bus consumes much energy. The cost of multiplying the two numbers is several orders of magnitude lower than that of accessing them from memory, creating a power wall.¹⁹ On the other hand, the difference in performance development between the two presents a storage wall. Processor performance increased by 50% per year compared to 7% per year for memory. The difference in data processing speed has become approximately 250 times greater than that in memory processing speed, as shown in Fig. 1.²⁰ The mismatch between the processor and memory results in slow

processing and wasted resources. In addition, the performance of a chip depends on the number of transistors. However, due to the physical limits of semiconductor processes, the size of the transistor reaches its limit. As the size of the device shrinks, the quantum and short-channel effects of the devices become increasingly severe. Overly dense transistors make heat dissipation more severe, and the power consumption of chips cannot be effectively reduced. These problems not only limit the performance of the chip but also have an impact on the accuracy and endurance of the calculations. People are desperate for new solutions.

The concept of computing-in-memory was developed in the last century, but the technology has not received much attention. The integration of computing power and storage units is costly and technically challenging. Integrated circuits are evolving at a rapid pace under Moore's Law. The needs of data processing are far from its limits. However, with the increased data volume and new demands for computing speed and accuracy, computing-in-memory technology is back in the



Fig. 1 Processor-memory performance gap.



Fig. 2 The development of the computing structure.²¹ (a) von Neumann architecture. (b) Near-memory computing structure. (c) Computing-in-memory structure.



spotlight. Computing-in-memory technology breaks through the traditional computing architecture, eliminating the need for a separate computing unit to process the data. The technology can be divided into charge-based and resistor-based storage technologies. Most of the charge-based storage technology is based on adding a computational circuit to a conventional memory cell, as shown in Fig. 2(b). For example, the SRAM memory cell is the central cell in a CMOS-based memory circuit. Simple logic circuits can be implemented by designing simple peripheral circuits, and efficient neural network training can also be achieved by SRAM across-array design. Resistor-based storage technologies use memory as the core of integrated circuits. The initial data for the system are the resistance values of the resistor cells, while the input data come from external electrical signals. The resistance value of each cell represents the final calculation result. Resistor-based circuits mainly use new non-volatile memories, such as RRAM, PCM, and FRAM. Advanced memory combined with an across-array structure enables complex computations such as vector-matrix multiplication, greatly accelerating computation speed and reducing operational power consumption. Resistor-based storage technologies will be described in detail below.

3 RRAM

3.1 Memristor

In 1971, Professor Chua^{10,22} of the University of California at Berkeley introduced the memristor concept and then extended its definition. He considered that a memristor is the fourth fundamental circuit element in addition to resistance, capacitance, and inductance. The memristor complements the gap in the fundamental circuit relationship between charge and flux. Its mathematical model is represented by the ratio of magnetic flux to charge. The memristor resistance is numerically equal to the ratio of the voltage applied to memory terminals at a given moment to the current flowing through it. In practice, however, the value of a memristor is determined by the necessary amount of current or flux flowing through the memristor over time, which is a non-linear resistive element with memory characteristics.²³

The prevailing explanation for the resistive principle of the memristor is the wire filamentary resistive switching mechanism. The wire filament mechanism assumes that the resistance change of the device corresponds to the formation and

rupture of wire filaments within the material.²⁴ In most cases, a freshly prepared memristor usually exhibits fewer defects and shows an initial resistance state (IRS), as shown in Fig. 3(a). When a positive voltage is applied to the device, the conductive defects in the resistive layer are connected to form conductive filaments. The device's conductivity is enhanced, showing a low resistance state (LRS), as shown in Fig. 3(b). When a negative voltage is applied to the device, the conductive filaments break, and the device's conductivity is reduced, resulting in a high resistance state (HRS), as shown in Fig. 3(c). Some of the memristors even have intermediate resistance levels or more resistance levels.

Researchers used the resistive properties of resistors to implement the calculation function. The memristor used in memory applications is called resistive random-access memory (RRAM). RRAM has high speed, high durability and multi-bit storage capability. It meets the development requirements of high capacity, high read and write times, faster read and write speeds, and lower power consumption. RRAM is expected to replace traditional storage devices to achieve a big boost in computing power.

3.2 Structure of RRAM

The basic structure of a memristor is a sandwich structure consisting of upper and lower metal electrodes and an intermediate insulator resistive layer. RRAM is constructed as an across-array structure to improve integration. Realization of high-density and reliable RRAM is crucial toward development of next-generation information storage and computing. The structure can be divided into active and passive arrays. In the active array, the RRAM forms an array with field effect tubes. The transistors control the reading, writing and erasing of resistive elements and can effectively isolate adjacent cells from interference. However, the active array design is complex and not conducive to integration. In passive arrays, each memory cell is defined by upper and lower electrodes formed by intersecting word and bit lines, as shown in Fig. 4(a). Passive arrays are used in most high-density RRAM arrays, and passive arrays are more conducive to the three-dimensional (3D) integration of the chip.

The across-array structure of RRAM is one of the bases of its fast calculations and responsible calculations. Employing a 3D structure can make resource-expensive tasks into a manageable size and provides substantial improvement to the speed and

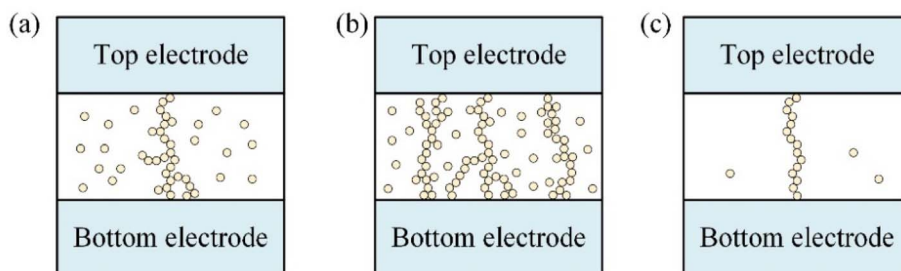


Fig. 3 The wire filament mechanism. (a) Initial resistance state (IRS). (b) Low resistance state (LRS). (c) High resistance state (HRS).





Fig. 4 RRAM cross-array architecture and scaling.²⁵ (a) RRAM cross-array structure with a single memory layer. (b) Horizontal stacked 3D cross-array structure. (c) Vertical 3D cross-array structure.

energy efficiency while running complex neural network models. A broad range of applications can be envisioned with further optimizations in the device performance and more carefully designed on-chip peripheral analogue circuitry, even though the current 3D design might not be the most scalable architecture compared to its two-dimensional (2D) counterpart. According to the stacking form, 3D RRAM can be divided into a horizontal stacked 3D cross-array structure (HRRAM) and vertical 3D cross-array structure (VRRAM). HRRAM is composed of multilayer planar two-dimensional RRAM

superimposed by the fabrication process, as shown in Fig. 4(b). VRRAM is made by rotating the conventional horizontal cross-array structure by 90° and extending it repeatedly in the horizontal direction, as shown in Fig. 4(c).

Three-dimensional fabrication technology not only increases the integration of the chip and greatly reduces the chip area, but also increases the computing speed and brings newer applications. However, making 3D chips is still a very difficult process, and only a few labs can make them. In 2022, Tang *et al.*²⁶ fabricated a RRAM cross-array using MoS₂ and then

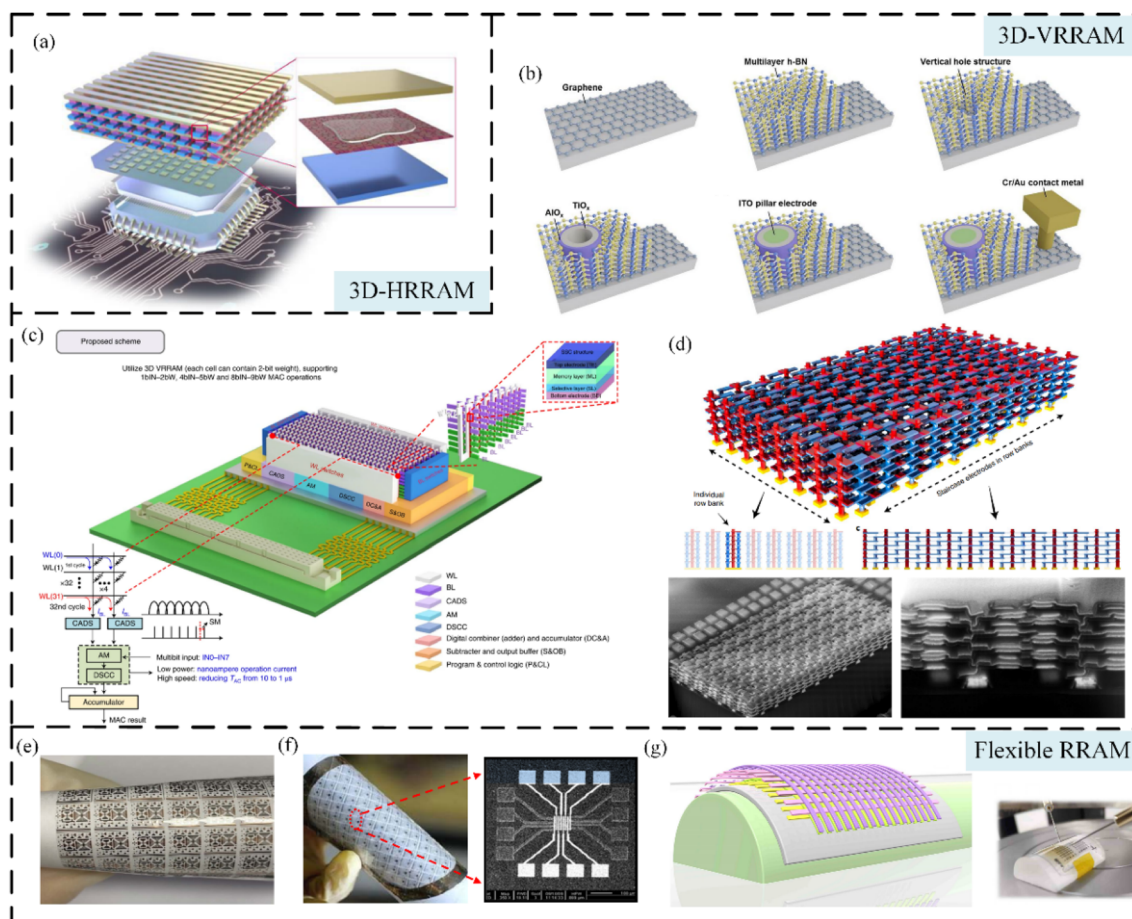


Fig. 5 Different forms of RRAM. (a) 3D-HRRAM: a 3D memristor array with buried metal interconnects and logic circuits. (b–d): 3D-VRRAM. (b) The fabrication flows of a single 2D material based vertical RRAM. (c) High-precision CIM scheme based on MLSS 3D VRRAM. (d) 3D monolithic integrated memristor circuits. (e and f): Flexible RRAM. (e) Flexible RRAM crossbar arrays with a Ti/PPX-C/Cu structure. (f) The flexible memory devices based on polychord-*para*-xylylene. (g) 3D flexible RRAM artificial synaptic network using a low-temperature atomic layer.



implemented a single three-dimensional RRAM cube by overlaying two-dimensional MoS₂ layers, as shown in Fig. 5(a). This work paves the way algorithmically for the implementation of two memristors in high density neuromorphic computing systems. In the process, it is possible to integrate the prepared 2D RRAM into 3D successfully. However, 3D-HRRAM requires critical lithography and other processes for each stacked layer, and this overhead manufacturing cost increases linearly with the number of stacks. The three-dimensional vertical structure occupies a smaller area in the integrated array of the same number of devices and is suitable for large-scale operations. Huang *et al.*²⁷ proposed a vertical architecture of RRAM design,

as shown in Fig. 5(b). The RRAM is composed of graphene plane electrode/multilayer hexagonal boron oxide (h-BN) insulating dielectric stacked layers, AlOx/TiOx, a resistive switching layer and an ITO pillar electrode that exhibits reliable device performance. The vertical three-dimensional structure combining the graphene plane electrode with a multilayer h-BN insulating dielectric can pave the way toward a new area of ultra-high-density memory integration in the future. Q. Huo *et al.*²⁸ reported a two-kilobit non-volatile CIM macro based on an eight-layer 3D VRRAM, as shown in Fig. 5(c). The chip is fabricated using a 55 nm complementary metal-oxide-semiconductor process. Scaling such systems to three-dimensional arrays

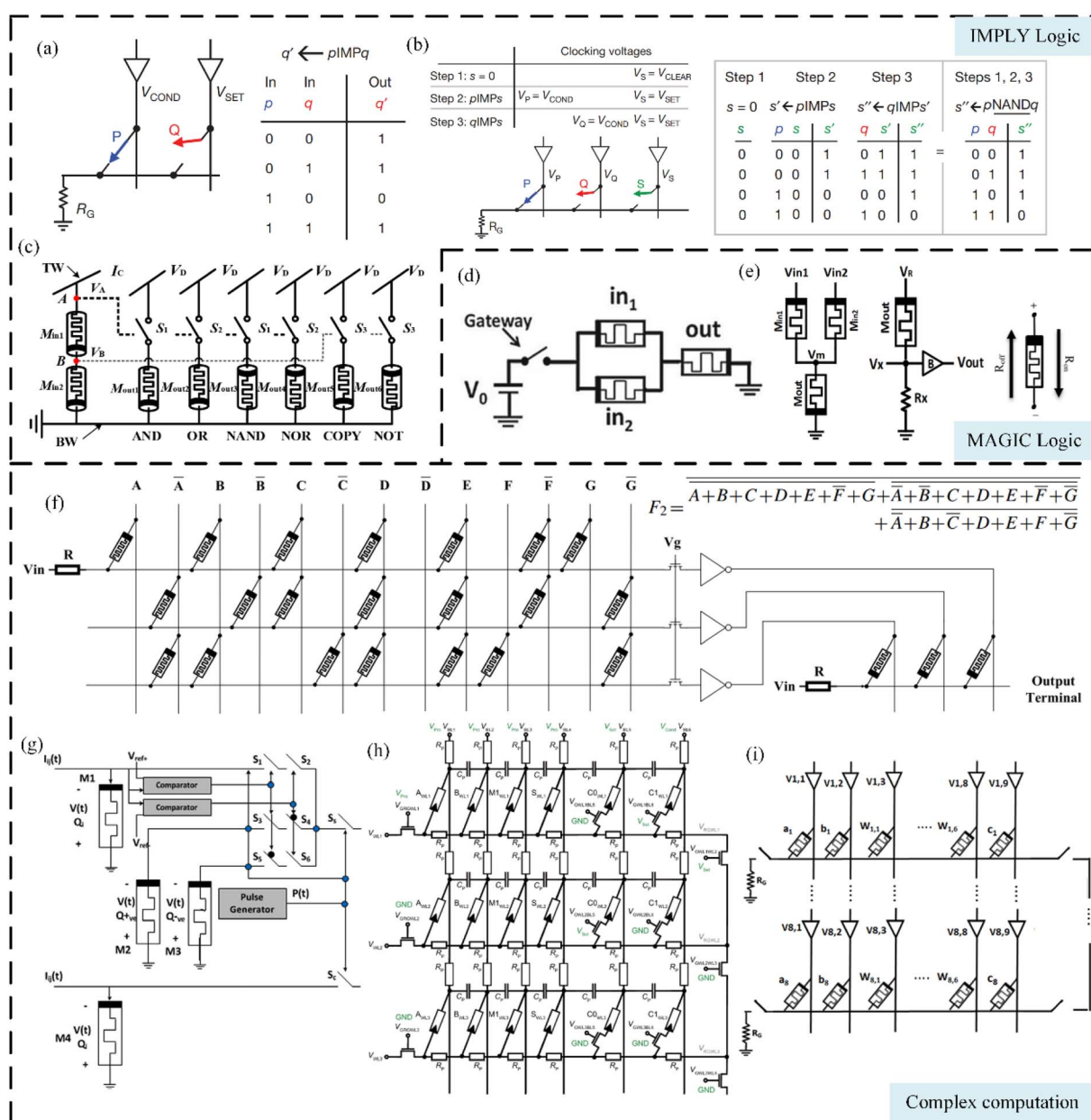


Fig. 6 Computing-in-memory realized by RRAM. (a) IMPLY logic and (b) NAND logic proposed by HP LABS. (c) The multi-functional boolean logic circuit schematic. (d and e): MAGIC logic computation. (d) MAGIC two-put OR gates and (e) NAND/NOR gates using the convert VTM method. (f–i): Complex computation. (f) Cross-array circuits for cases beyond the constraints. (g) Memristor-based quinary half adder. (h) Two-bit adder circuit with parasitic elements based on RRAM. (i) Digital full adder based on RRAM.



could provide higher parallelism, capacity and density for necessary vector–matrix multiplication operations. Lin *et al.*²⁹ successfully constructed an eight-layer array, as shown in Fig. 6(d). The 3D array was fabricated on top of a silicon wafer with a 100 nm thick thermally grown oxide layer. Researchers programmed parallelly operated kernels into the prepared RRAM, implemented a convolutional neural network and achieved software-comparable accuracy in recognizing handwritten digits.

In addition, RRAM, made of flexible materials for particular application scenarios, is also a hot topic. Kim *et al.*³⁰ used PPXC as a resistive switching layer and substrate to fabricate RRAM flexible across-arrays, as shown in Fig. 5(e). Researchers applied them to neuromorphic accelerated computing to test their performance with high electrical conductivity and durability. Huang *et al.*³¹ proposed a kind of single-component polymer RRAM based on polychord-*para*-xylylene, as shown in Fig. 6(f). The device has excellent chemical stability and high CMOS process compatibility as well as further reduced operation current. Wang *et al.*³² successfully prepared a 3-layer flexible RRAM array by a low-temperature atomic layer deposition technique to achieve high-density binary storage function and multi-bit storage in a single device, as shown in Fig. 5(g). Preparing flexible devices opens up new possibilities for future high-density, high-performance wearable applications.

A RRAM across-array can greatly compress the chip volume and improve the computing efficiency and computing volume. But the problem is also obvious: chip design and preparation processes are difficult, every application needs to be customized, there is still a lot of room for development and there are still unresolved problems.

(a) When writing and programming, current inevitably passes through other resistive cells. Leakage current can interfere with the accuracy of the information read when adding unnecessary power consumption.

(b) Due to the fabrication process, the ion movement pattern inside the device is somewhat random and not perfectly uniform for each resist-variable cell.³³

(c) When the device is operating, thermal crosstalk affects the change in the resistance state of the surrounding resistive cells—leading to inaccurate calculations. The effect of thermal crosstalk even increases with increasing density.

4 Applications

RRAM is considered one of the outstanding candidates for emerging storage technologies. It has low power consumption, high speed computing, high density and non-volatile data storage capability. RRAM has been applied in logic computing, neural networks, brain-like computing, and fused technology of sense-storage-computing.

4.1 Computing-in-memory

The simplest application is logic operations; the current application of logic circuits implemented logic gate operations using resistive cells instead of traditional transistors. HP LABS³⁴ first

prepared a physical model of RRAM in 2008 and implemented IMPLY logic operations. The primary logic cell and truth table are shown in Fig. 6(a). IMPLY and NOT logic can form a complete logic set to complete all 16 logics. However, it is not a beautiful operation, which destructs the input method because the final output of the implied logic overwrites the original input value. Implementing other logic using IMPLY logic also requires more steps, which is less efficient. Then, HP LABS also proposed the 3M1R structure, as shown in Fig. 6(b). It is an improvement to IMPLY, whose basic logic unit consists of three resistive cells and a voltage divider resistor. The high resistance state represents logic “1” and the low resistance state represents logic “0”, where P and Q represent the input variables and S represents the output variables. Based on the research from HP LABS, more complex logic operations can be obtained using RRAM. Dong *et al.*³⁵ combined multiple logic resistors in a single circuit to implement logic operations, as shown in Fig. 6(c), which offers the possibility of implementing advanced computing architectures. The study implements two-input or multi-input AND, OR, NOR, and NAND operations and single-input copy and NOT operations. In each logic gate, the circuit uses non-volatile resistors of the RRAM as input and output states, thus enabling stateful logic operations. In contrast to several existing methods, this method generates a versatile state logic circuit that can perform multiple state logic operations simultaneously.

In addition to the structure proposed by HP LABS, researchers connected RRAM in series and parallel. S. Kvatinsky *et al.*³⁶ proposed a memristor-aided logic (MAGIC) circuit that uses resistive values as logical state variables, as shown in Fig. 6(d). Under applied voltage control, the five essential logic functions can be achieved by connecting the RRAM in series and parallel. This design approach allows the results of logic calculations to be stored in a separate resistive cell, avoiding the problem of data overwriting arising in IMPLY logic operation. F. Mozafari *et al.*³⁷ proposed a new memristor-based NAND/NOR logic gate with a similar structure for general-purpose logic gates requiring two different input voltages. The structure consists of two input amnesia resistors and one output amnesia resistor, as shown in Fig. 6(e). RRAM-based logic calculations can reduce the computational complexity of the problem, as well as reduce the amount of data being accessed by performing the calculations within the across-array.

In recent years, the improvement of the technological level and the progress of materials research have enabled the successful preparation of larger RRAM matrices and the realization of more and more complex functions. Cui *et al.*³⁸ proposed a family of NMOS-like RRAM gates. All gates are array implementable, and the gate family logic is complete. NOR, AND and NOR gates consume only one cycle during the computation phase. The gate circuit saves half the number of RRAM devices compared to its CMOS-based counterpart. By using RRAM gates and across-array structures, complex logic operations can be realized, as shown in Fig. 6(f). A. H. Fouad *et al.*³⁹ proposed a multi-valued logic adder, as shown in Fig. 6(g). The study discussed the possibility of extending a three-valued adder circuit to a multi-valued logic adder



theoretically. By exploiting the properties and dynamics of RRAM, the circuit achieved the advantages of handling different numbering systems, increasing density, and reducing processing time. Simon *et al.*⁴⁰ proposed a multi-input memristive switch logic, as shown in Fig. 6(h). This work enabled the function XOR (Y NOR Z) to be performed in a single step with three memristive switches. This OR/NOR logic gate increases the capabilities of memristive switches, improving the overall system efficiency of a memristive switch-based computing architecture. N. Taherinejad *et al.*⁴¹ present a new architecture for a digital full adder, as shown in Fig. 6(i). The circuit is faster than existing IMPLY-based serial designs while requiring less area compared to the existing parallel design.

In the latest studies, researchers hope not only to implement logical operations, but also to implement complex mathematical analysis in RRAM cross-arrays. Tian *et al.*⁴² realized a hardware Markov chain algorithm in a single device for machine learning. M. Teimoori *et al.*⁴³ designed a 2×2 multiplier circuit, as shown in Fig. 7(a). The circuit requires only sixteen resistive cells, eight transistors, and only one calculation time step for the multiplication operation, which is of low cost. B. Chakrabarti *et al.*⁴⁴ demonstrated vertical monolithic

integration of 3D RRAM crossbars on a foundry-processed 5×5 CMOS chip. The chip can realize dot-product operation, which is the first application of functional 3D CMOL hybrid circuits. Xie *et al.*⁴⁵ used h-BN RRAM cross-arrays to demonstrate the hardware implementation of dot product operation and the linear regression algorithm, as shown in Fig. 7(b)–(d). Y. Hala-wani *et al.*⁴⁶ proposed an XNOR-based RRAM content addressable memory (CAM) with an analog time-domain adder function for efficient winning class computation, as shown in Fig. 7(e). The chip had 31 times less area and about three times less energy consumption than state-of-the-art RRAM designs. P. Kumar *et al.*⁴⁷ used h-BN as an electrolyte to implement the fabrication of an across-array RRAM, as shown in Fig. 7(f), and combined it with CMOS circuits to implement extreme learning machine algorithms. With CMOS, circuits implement the encoder unit and RRAM arrays implement the decoder unit. The hybrid architecture performs well on complex audio, image and non-linear classification tasks with real-time data sets.

These studies are of great significance for further implementation of neuromorphic computing and machine learning hardware based on RRAM cross-arrays. The compatibility of the above design with CMOS technology provides a new way for

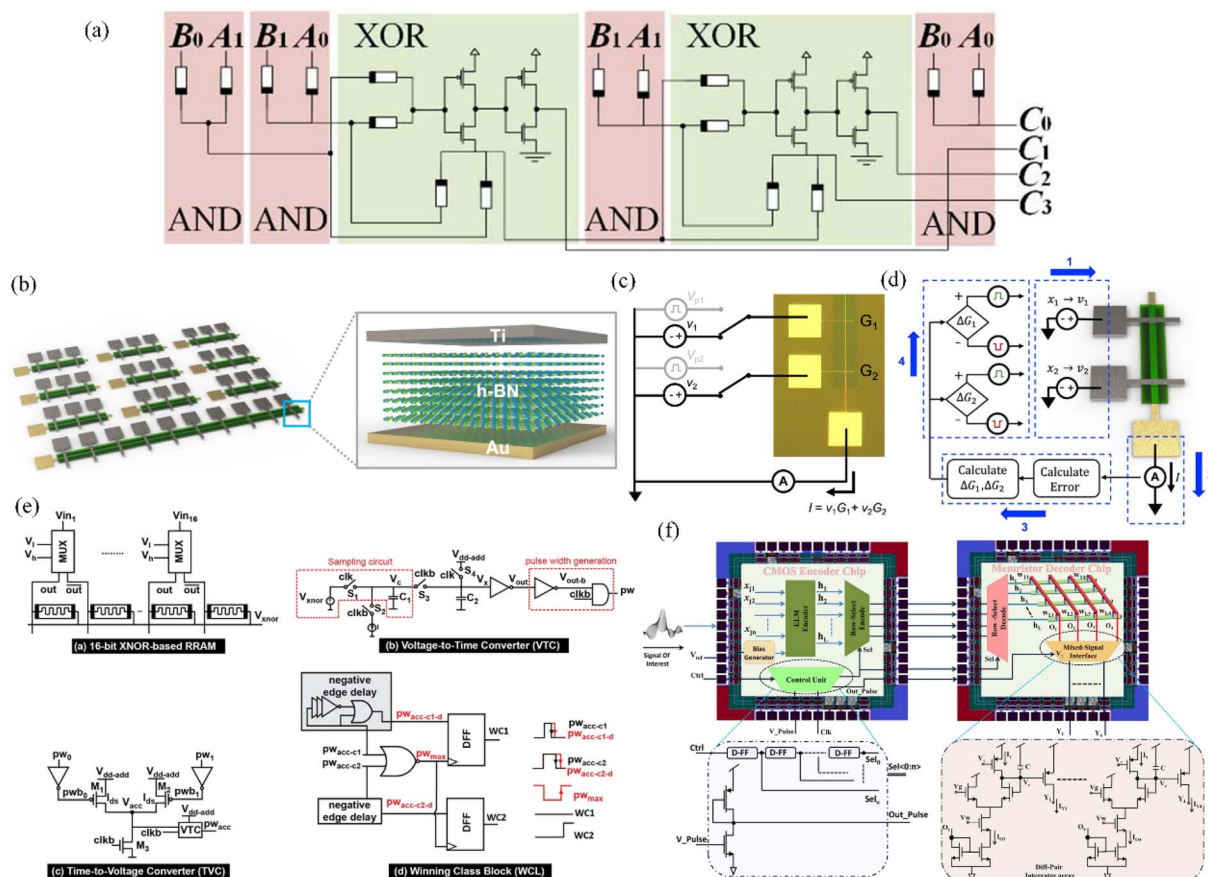


Fig. 7 Complex mathematical formulas realized by RRAM. (a) Logic design of the 2×2 multiplier. (b) Schematic of Au/h-BN/Ti RRAM arrays and cross-sectional schematic of a single memristor. (c) Schematic of the experimental setup for demonstration of basic dot-product operation on h-BN RRAM arrays. (d) Flow diagram for stochastic multivariable linear regression. (e) Circuit designs of the proposed RRAM-based CAM, analog time-domain adder and winning class logic. (f) System architecture showing a CMOS encoder chip and memristor decoder chip along with its functional sub-blocks.



the development of digital logic circuits. Introduction of RRAM was presented as a solution to the fundamental limitations of VLSI systems.

4.2 Neural networks and brain-like computing

In recent years, RRAM has been widely used in neural morphological networks. Artificial neural networks and all their variants are the primary tools for machine learning tasks. The continuous production of large amounts of data makes it possible to train and operate artificial neural networks successfully. However, in order to achieve high inference accuracy, neural networks usually require a large number of parameters. The traditional structure of memory and computing points makes the whole stage consume a lot of time and energy but using RRAM to build neural networks is expected to solve this problem. RRAM has resistive properties and is similar to synapses in neural networks, so the neural network hardware and a neural network model based on RRAM have intrinsic consistency. A single device can be used as a synapse and can integrate storage and operation. A RRAM-based neural network integrates computing and storage closely, eliminating

data transmission between the processor and memory, thus improving the overall system performance and saving most system energy consumption. RRAM is well compatible with CMOS processes and can be scaled up on a large scale through cross-array structures. Through CMOS processes, functional circuits are added to the periphery of RRAM, so that the system can complete the calculation of a matrix in one operation, thus playing an important role in high dimensional computation.

The perceptron model is an early neural network algorithm. In the calculation process of the perceptron neural network, matrix vector multiplication between the input information vector and weight matrix consumes a lot of computing resources. However, using RRAM cross-arrays to realize matrix vector multiplication in parallel in one step can greatly reduce the energy consumption of the hardware neural network. The RRAM cross-array is used to store the synaptic weight matrix, and the conductance at each crossing is used to represent the weight value of a synaptic connection. MNIST image classification is modeled using a single-layer perceptron that is not *in situ* trained and a large-scale multi-layer perceptron classifier based on a more advanced conductance tuning algorithm. M.

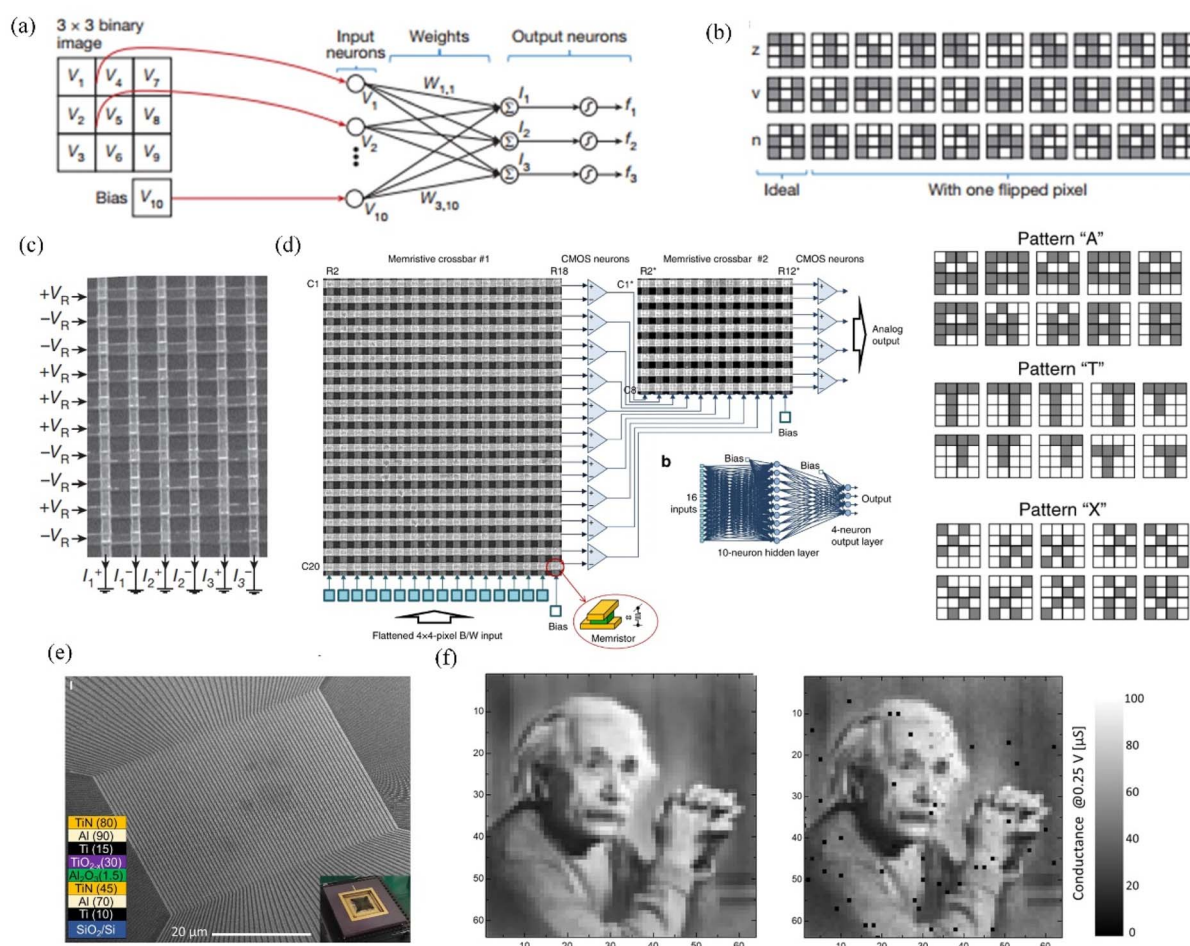


Fig. 8 Perceptron network based on RRAM. (a–c): Single-layer perceptron network. (a) Top-level description of the network. (b) The used input pattern set. (c) 12×12 RRAM cross-array circuit. (d–f): Multilayer perceptron network. (d) Double layer 20×20 RRAM cross-array circuit combined with the CMOS process. (e) 64×64 RRAM passive cross-array circuit. (f) Programming Einstein image to (e) with a 5% relative error.



Prezioso *et al.*⁴⁸ first produced an across-array of RRAMs for neural network computing. Researchers used a 12×12 RRAM across-array to form a single-layer perceptron network to recognize and classify black and white letters in 3×3 pixels, as shown in Fig. 8(a). F. M. Bayat *et al.*⁴⁹ produced two 20×20 RRAM across-array, as shown in Fig. 8(d), and integrated them with discrete CMOS components on two printed circuit boards to implement a multilayer perceptron (MLP) neural network. The MLP network has 16 inputs, 10 hidden layer neurons, and 4 outputs, enough to classify 4×4 pixel black and white patterns into 4 categories. Then, more complex array structures were proposed based on previous research to realize much more complex computing. H. Kim *et al.*⁵⁰ developed a 64×64 passive crossbar circuits with almost 99% working cross point metal-oxide memristors based on the foundry-compatible fabrication process. The circuit has a relatively average error of <4% when programming 4 K grey-scale patterns and an error of nearly 1%. It is an essential step towards realizing a human brain-scale integrated neuromorphic system.

The convolutional neural network (CNN) reduces the dimensions by convolutional operation, reduces parameter complexity, simplifies complex problems before processing, and greatly reduces the calculation cost. The convolutional neural network based on memristor cross arrays mainly consists of two parts: the convolutional operation part and the fully connected layer part. On the one hand, the memristor array can store convolution kernels and complete the matrix vector multiplication of input information and convolution kernels in one step, which greatly improves the computational efficiency. On the other hand, the fully connected layer part of a convolutional neural network is a multi-layer perceptron, which can also be realized in parallel by using the memristor cross array as mentioned above. However, CNNs have not yet fully implemented hardware *via* a RRAM across-array. P. Yao *et al.*⁵¹ implemented complete hardware for convolutional networks based on RRAM array chips. The team used eight arrays with 2048 RRAM to implement a five-layer convolutional neural network. The chip was two orders of magnitude more energy efficient than an image processor (GPU) in processing convolutional neural networks by comparison with conventional neural network computation.

Stochastic neural networks introduce random changes into neural networks, one is to assign transfer functions of the stochastic process among neurons and the other is to assign random weights to neurons. This makes stochastic neural networks very useful in optimization problems, because random transformations avoid local optimality. The key operation in stochastic neural networks is the random dot product, which has become the latest method to solve the problems of machine learning, information theory and statistics. Although there have been many demonstrations of dot product circuits and random neurons, there is still a lack of efficient hardware implementation combining these two functions. M. R. Mahmoodi *et al.*⁵² proposed versatile stochastic dot product circuits based on nonvolatile memories for high performance neuro-computing and neuro-optimization, as shown in Fig. 9. F. Zahari *et al.*⁵³ proposed an analogue pattern recognition with stochastic switching binary CMOS-integrated RRAM. Researchers prepared a random binary RRAM device based on a polycrystalline HfO and have produced software neurons. Based on this, a random neural network is used for image recognition. The convergence rate of this new learning algorithm is very fast, which may significantly reduce the power consumption of the system. This study has potential applications in stochastic neural networks to solve MNIST pattern recognition tasks.

Reservoir computing (RC) is a concept developed from recursive neural networks (RNNs) that has recently been successfully used to implement a wide range of tasks, such as image model recognition, time series prediction, and pattern generation. Researchers have successfully used memristors for this purpose as well. Fig. 10 shows the concept of a memristor-based RC system. The spikes collected from firing neurons are used directly as inputs to an excitation memristor. The repository space is further extended with the concept of virtual nodes to help handle complex time inputs. A simple neural network is used as the reservoir's readout layer to produce the final output. In 2020, Lu *et al.*⁵⁴ looked at reserve pool arithmetic. Researchers demonstrate a RC system based on RRAM, whose states reflect the temporal characteristics of nerve peak sequences, successfully used to identify neural firing patterns, monitor the conversion of firing patterns, and identify neural



Fig. 9 Stochastic neural networks.⁵² (a) Stochastic dot-product circuit and its applications in neurocomputing. (b) RRAM-based restricted Boltzmann machine. (c) A bipartite graph of the considered RBM network and its implementation (the red rectangle highlights the utilized area of a 20×20 across-array).





Fig. 10 Schematic showing the concept of a memristor-based RC system for neural activity analysis.

synchronization states between different neurons. This work makes it possible to realize efficient neural network signal analysis with high spatial and temporal accuracy. In addition, researchers made a RRAM across-array based on WO_x . They used a circuit to form an array to recognize digits and increased the noise, proving the anti-interference and stability of the system's recognition, and expanded to a small reserve pool composed of 88 memristors for handwritten digit recognition. X. Liang *et al.*⁵⁵ also applied RRAM to reservoir computing, an artificial neural network for efficient processing of temporal signals, with the potential to have a significant impact in the field of brain-like computing.

On the basis of the development of neural networks, the concept of brain-like computing has been put forward. From an

information technology perspective, researchers have abstracted the way decisions are made in the human brain. They connected it through multi-layered artificial neural networks to create a non-linear, adaptive computing model for information processing. Wu *et al.*⁵⁶ have been focusing on research on brain-like computing for a long time and have described the RRAM-based brain-like circuit in detail, as shown in Fig. 11(a). By introducing artificial dendritic computational units with rich dynamic properties, the team constructed a novel artificial neural network,⁵⁷ significantly reducing system power consumption while improving computational network accuracy. Compared to traditional processes, the new system has a 30 times reduction in dynamic power consumption, an 8% increase in accuracy, and an overall system power

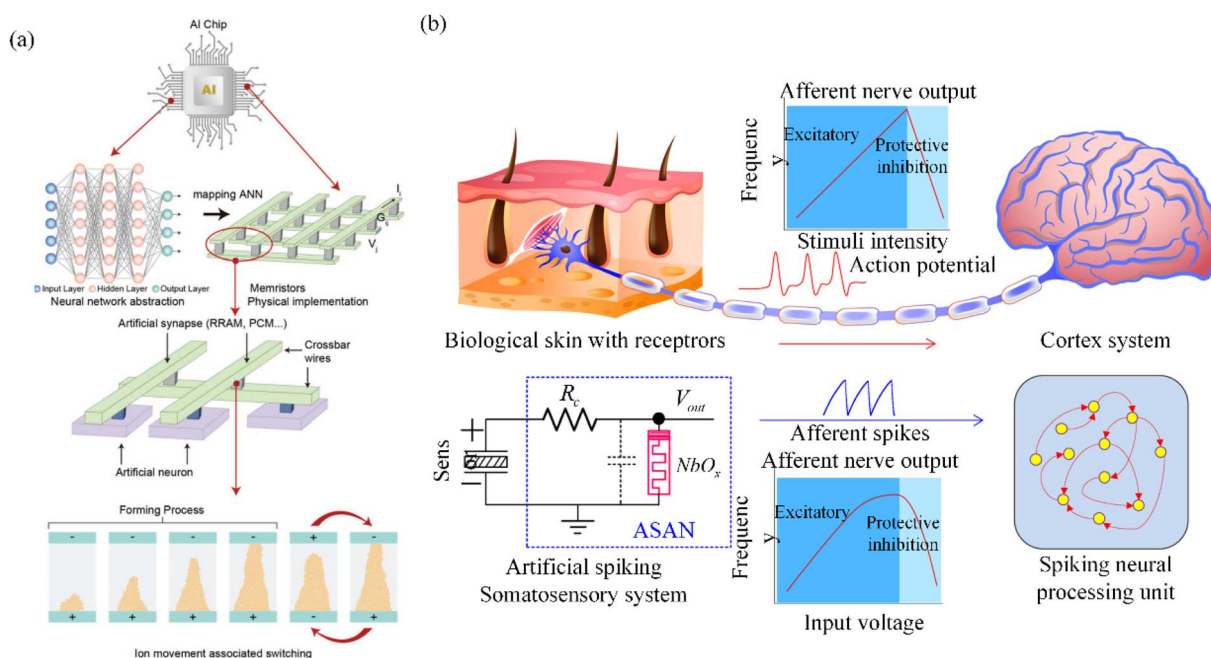


Fig. 11 Brain-like computing. (a) Neural network operations using RRAM implementation. (b) Biological afferent nerve vs. the artificial afferent nerve.



consumption downgraded by three orders of magnitude. Z. Wu *et al.*⁵⁸ built an artificial sensory neural system implementation scheme with habituation properties based on RRAM, as shown in Fig. 11(b). The system used habituation as a biological learning rule to build a habituated impulse neural network that can be applied to robot autonomous cruising obstacle avoidance. The neural network can be used to implement robot obstacle avoidance functions and effectively improve robot obstacle avoidance efficiency.

The application of RRAM-based neural networks is an indispensable research direction in neuromorphic computing. As a new type of synaptic device, the memristor realizes the hardware mapping of the synaptic weight well by means of multivalued modulation under pulse and can store the synaptic weight matrix and realize *in situ* calculations. The parallel matrix multiplication and addition capability of the across-array structure of RRAM realizes the acceleration of neural network computing and provides a solution for constructing a new computing architecture.

4.3 Fused technology of sense-storage-computing

The human nervous system senses the physical world in an analogue but efficient way. Researchers have been working on using hardware to simulate human perceptual systems. Sensors take signals such as light, pressure and sound and convert them into electrical signals. The processor is used to analyze electrical signals and simulate human sensory systems. However, traditional computing platform analysis speed is slow and precision is not high. As a new type of memory, RRAM can be made of materials that are sensitive to peripheral signals and can itself be used as a sensor. In addition, RRAM can combine micro-sensors with neural network computing. The sense-storage-computing fusion technology based on RRAM brings a new research direction for bionics technology.

Visual perception is made by using the effect of light signals on RRAM or using the RRAM across-array as a processing circuit, which can currently be used to identify numbers and test light intensity. Emboras *et al.*⁵⁹ first discovered the change



Fig. 12 Fused technology of sense-storage-computing. (a) Schematic structure of ORRAM. (b) Schematic structure of an 8×8 ORRAM array. (c) Three-dimensional schematic of the investigated optically readable plasmonic RRAM.⁵⁹ (d) The common implementation scheme for sound localization application with CMOS circuits. (e) A conceptual diagram of a memristor-based neuromorphic sound localization system. The RRAM across-array acts as synapses to deal with binaural sound signals. (f) The conceptual diagram of haptic memory. (g) Image for 7×7 SS-RRAM arrays and ECG electrode mounted onto the wrist. (h) Structure of a memristor for gas sensing applications.



in light signal transmission in the waveguide during transformation, which laid the foundation for visual perception based on RRAM. Zhou *et al.*⁶⁰ prepared a photo-resistive memory, as shown in Fig. 12(a) and (b). The component can convert light signals into electrical signals and store the data briefly. In conjunction with developments in neural networks, researchers used RRAM as a synaptic node to enable more advanced visual processing. Lorenzi *et al.*⁶¹ used an RRAM array to recognize a 5×5 pixel binary picture. Sarkar *et al.*⁶² used an RRAM-based neuromorphic circuit system and trained and tested the system to successfully recognize basic Arabic numerals. Not only digital numbers, but also visual recognition of infrared, low light and images are currently being investigated.

In terms of sound localization, researchers have localized the location of sound sources based on the theory of interaural time difference, where the time difference between the theoretical sound arriving in the two ears is shown in Fig. 12(d). Gao *et al.*⁶³ have a leading study in this area. The team fabricated a 128×8 array of memristors and proposed a new sound localization algorithm. The devices were tested with pulses to simulate sound signals, showing that the training accuracy and energy consumption of the system built using the memristor array were significantly improved, representing a significant advance in auditory localization systems with memristors. F. Moro *et al.*⁶⁴ designed an event-driven object localization system inspired by the owl auditory cortex using a combination of state-of-the-art piezoelectric micromechanical ultrasound sensors and RRAM. The system was more efficient and power-efficient than microcontrollers performing the same task by several orders of magnitude.

Haptic perception is mainly studied in terms of sensing external pressure, where pressure data are converted into electrical signals and stored in memory. The location of pressure generation can be clarified through an array mechanism of resistive variable memory, as shown in Fig. 12(f). Park *et al.*⁶⁵ used a flexible RRAM array for ECG measurements, as shown in Fig. 12(g). Researchers successfully demonstrated stable data storage of cardiac signals, a damage-reliable memory triggering system using a triboelectric energy-harvesting device, and touch sensing *via* pressure-induced resistive switching.

As for smelling perception, the research started late, and the system can identify a small type of gas. People use the combination of a gas sensor and RRAM made of special materials and use the storage array to store and process the data quickly, so as to achieve the purpose of recognition. A. Adeyemo *et al.*⁶⁶ used an across-array RRAM as a gas sensor, as shown in Fig. 12(h). Researchers used the HP LABS fabricated TiO₂ based memristor model in an attempt to improve sensing accuracy. The experiment provides the basis for the research of smelling perception. At present, the smell sensing system using RRAM has a long sensing period, low sensing accuracy, single sensing gas, and narrow application. Most systems have a small response gap when identifying similar gases, which requires additional amplifying circuits to be sufficient for classification, increasing the structural complexity.

4.4 Other applications

RRAM has good switching and isolation characteristics for high-frequency communications. In 5G communication, the RRAM has 50 times the switching efficiency of other non-volatile switches. At 6G operating frequencies, RRAM also exhibits high isolation and has sub-nanosecond pulse switching, low BER and a high signal-to-noise ratio. M. Kim *et al.*^{67–69} worked on RRAM for high-frequency applications at the University of Texas. RRAM made from MoS₂ and h-BN is suitable as an energy-efficient RF switch, overcoming the limitations of transistors and pick-and-place switches. Researchers also predicted that RRAM could be used in mobile systems, low-power IoT and terahertz beam steering. The unique switching randomness of RRAM has important applications in information security and can also be used to design chaotic circuits. H. M. Ibrahim *et al.*⁷⁰ have implemented a robust and lightweight physically unclonable function (PUF) architecture using RRAM, with implications for cryptographic key and security application upgrades. In addition, P. S. Zarrin *et al.*⁷¹ proposed a neuromorphic on-chip recognition of saliva samples of COPD and healthy controls using RRAM. This is an application of RRAM in the medical field.

5 Conclusion and outlook

With the constraints of the von Neumann structure and the physical limits of semiconductor processing, computers are experiencing bottlenecks in their ability to process data. As a new computing principle, computing-in-memory technology embeds the computation capability into the storage unit. This increases computation speed, reduces system power consumption, and opens up more application possibilities.

RRAM is a new memory type with resistive and non-volatile properties. The structure, mechanism and preparation of RRAM have been the focus of attention, with research continuing in in-memory logic operations, brain-like computing and fused technology of sense-storage-computing. However, research into RRAM still faces some challenges:

(a) RRAM has randomness. The characteristics presented by memory resistors are not identical from cycle to cycle. There is also variability in devices made from the same batch. The preparation process of RRAM is being improved to try to solve this problem. In addition, the addition of peripheral circuits to the memristor array, the use of check methods and redundant design are also used to reduce errors.

(b) The best materials for the resistive layer in RRAM are also being screened. Metal oxides, 2d materials, emerging materials and organics are currently used for the preparation of RRAM. The materials used to prepare devices largely influence the performance of amnesic resistors. High resistance ratios, good homogeneity and matching proven manufacturing processes and equipment are the basis for judging the suitability of the material. We need to develop new materials with innovative approaches and explore the best material composition ratios based on already well-defined material elements. Better materials will enhance the performance given to each application.



(c) Mass production is almost impossible to replace most conventional memories. Across-array circuits are still in the small-scale laboratory stage. The preparation of memristors is difficult, with only a few laboratories and processors meeting the demand. The variety of materials used in the preparation of memristors and the preparation processes involved each have their own characteristics, and integration with existing mature processes takes time. Complex functional applications require customization. Different functions can be achieved by selecting different memory resistor units.

Data availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

Conflicts of interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This research was supported by the National Natural Science Foundation of China under Grant 61971389, Grant 62131010, and Grant 61801253 and in part by the Zhejiang Provincial Natural Science Foundation of China under Grant LR22F010001, Grant LZ20F010006, Grant LY20F010003 and LY21F010006 and the Fundamental Research Funds for the Provincial Universities of Zhejiang Grant 2020YW01.

References

- M. Lundstrom, Moore's Law Forever, *Science*, 2003, **299**(5604), 210–211.
- R. R. Schaller, Moore's law: past, present and future, *IEEE Spectrum*, 1997, **34**(6), 52–59.
- J. Sim, M. Kim, Y. Kim, *et al.*, *MAPIM: Mat Parallelism for High Performance Processing in Non-volatile Memory Architecture; Proceedings of the 20th International Symposium on Quality Electronic Design (ISQED)*, F 6-7 March 2019, 2019.
- W. H. Kautz, Cellular Logic-in-Memory Arrays, *IEEE Trans. Comput.*, 1969, **C-18**(8), 719–727.
- A. Jaiswal, I. Chakraborty, A. Agrawal, *et al.*, 8T SRAM Cell as a Multibit Dot-Product Engine for Beyond von Neumann Computing, *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2019, **27**(11), 2556–2567.
- R. Guo, L. You, Y. Zhou, *et al.*, Non-volatile memory based on the ferroelectric photovoltaic effect, *Nat. Commun.*, 2013, **4**(1), 1990.
- M. Wuttig and N. Yamada, Phase-change materials for rewriteable data storage, *Nat. Mater.*, 2007, **6**(11), 824–832.
- W. Li, X. Qian and J. Li, Phase transitions in 2D materials [J], *Nat. Rev. Mater.*, 2021, **6**(9), 829–846.
- C. Chappert, A. Fert and F. N. Van Dau, The emergence of spin electronics in data storage, *Nat. Mater.*, 2007, **6**(11), 813–823.
- L. Chua, Memristor-The missing circuit element [J], *IEEE Trans. Circuit Theory*, 1971, **18**(5), 507–519.
- D. Strukov, G. Snider, D. Stewart, *et al.*, The Missing Memristor Found, *Nature*, 2008, **453**, 80–83.
- I. Hossen, M. A. Anders, L. Wang, *et al.*, Data-driven RRAM device models using Kriging interpolation, *Sci. Rep.*, 2022, **12**(1), 5963.
- J. Kang, T. Kim, S. Hu, *et al.*, Cluster-type analogue memristor by engineering redox dynamics for high-performance neuromorphic computing, *Nat. Commun.*, 2022, **13**(1), 4040.
- D. Ielmini, Resistive switching memories based on metal oxides: mechanisms, reliability and scaling, *Semicond. Sci. Technol.*, 2016, **31**(6), 063002.
- S. Thomas, IGZO and RRAM team up, *Nat. Electron.*, 2020, **3**(7), 353.
- E. C. Ahn, H. S. P. Wong and E. Pop, Carbon nanomaterials for non-volatile memories, *Nat. Rev. Mater.*, 2018, **3**(3), 18009.
- S. Bulja, R. Kopf, A. Tate, *et al.*, High frequency resistive switching behavior of amorphous TiO₂ and NiO, *Sci. Rep.*, 2022, **12**(1), 13804.
- H. S. P. Wong and S. Salahuddin, Memory leads the way to better computing, *Nat. Nanotechnol.*, 2015, **10**(3), 191–194.
- M. Horowitz. *1.1 Computing's Energy Problem (And what We Can Do about it) [M]*. 2014.
- D. Patterson, T. Anderson, N. Cardwell, *et al.*, A case for intelligent RAM, *IEEE Micro*, 1997, **17**(2), 34–44.
- A. Sebastian, M. Le Gallo, R. Khaddam-Aljameh, *et al.*, Memory devices and applications for in-memory computing, *Nat. Nanotechnol.*, 2020, **15**(7), 529–544.
- L. Chua, If it's pinched it's a memristor, *Semicond. Sci. Technol.*, 2014, **29**(10), 104001.
- L. Chua. Everything You Wish to Know About Memristors but Are Afraid to Ask [M]/CHUA L, SIRAKOULIS G C, ADAMATZKY A, *Handbook of Memristor Networks*, Springer International Publishing, Cham, 2019, pp. 89–157.
- Y. Li, S. Long, Y. Liu, *et al.*, Conductance Quantization in Resistive Random Access Memory, *Nanoscale Res. Lett.*, 2015, **10**(1), 420.
- D. Ielmini and H. S. P. Wong, In-memory computing with resistive switching devices, *Nat. Electron.*, 2018, **1**, 333–343.
- B. Tang, H. Veluri, Y. Li, *et al.*, Wafer-scale solution-processed 2D material analog resistive memory array for memory-based computing, *Nat. Commun.*, 2022, **13**(1), 3037.
- Y.-J. Huang and S.-C. Lee, Graphene/h-BN Heterostructures for Vertical Architecture of RRAM Design, *Sci. Rep.*, 2017, **7**(1), 9679.
- Q. Huo, Y. Yang, Y. Wang, *et al.*, A computing-in-memory macro based on three-dimensional resistive random-access memory, *Nat. Electron.*, 2022, **5**(7), 469–477.
- P. Lin, C. Li, Z. Wang, *et al.*, Three-dimensional memristor circuits as complex neural networks, *Nat. Electron.*, 2020, **3**(4), 225–232.



- 64 F. Moro, E. Hardy, B. Fain, *et al.*, Neuromorphic object localization using resistive memories and ultrasonic transducers, *Nat. Commun.*, 2022, **13**(1), 3506.
- 65 J. Park, D. Seong, Y. J. Park, *et al.*, Reversible electrical percolation in a stretchable and self-healable silver-gradient nanocomposite bilayer, *Nat. Commun.*, 2022, **13**(1), 5233.
- 66 A. Adeyemo, A. Jabir, J. Mathew, *et al.* *Reliable Gas Sensing with Memristive Array [M]*. 2017.
- 67 M. Kim, R. Ge, X. Wu, *et al.*, Zero-static power radio-frequency switches based on MoS₂ atomristors, *Nat. Commun.*, 2018, **9**(1), 2524.
- 68 M. Kim, E. Pallecchi, R. Ge, *et al.*, Analogue switches made from boron nitride monolayers for application in 5G and terahertz communication systems, *Nat. Electron.*, 2020, **3**(8), 479–485.
- 69 M. Kim, G. Ducournau, S. Skrzypczak, *et al.*, Monolayer molybdenum disulfide switches for 6G communication systems, *Nat. Electron.*, 2022, **5**(6), 367–373.
- 70 H. M. Ibrahim, H. Abunahla, B. Mohammad, *et al.*, Memristor-based PUF for lightweight cryptographic randomness, *Sci. Rep.*, 2022, **12**(1), 8633.
- 71 P. S. Zarrin, F. Zahari, M. K. Mahadevaiah, *et al.*, Neuromorphic on-chip recognition of saliva samples of COPD and healthy controls using memristive devices, *Sci. Rep.*, 2020, **10**(1), 19742.

