

# In-Memory Computing-Based Ultra-Efficient Massive MIMO Precoding: Memristor Crossbar Circuits, Conductance Mapping Strategies, and Programming Latency Estimation

Yu-Xin Zhang<sup>1</sup>, Shaoshi Yang<sup>1</sup>, *Senior Member, IEEE*, Yi-Hang Ren, Sheng Chen<sup>2</sup>, *Life Fellow, IEEE*, and Ping Zhang<sup>1</sup>, *Fellow, IEEE*

**Abstract**—Matrix inversion in massive multiple-input multiple-output (MIMO) precoding imposes significant burden on the energy efficiency and processing latency of baseband units. In this paper, we first propose a memristor crossbar based in-memory computing circuit capable of supporting both zero-forcing (ZF) and minimum mean square error (MMSE) precoding. The circuit features a reduced matrix size and enables faster one-step computation without the need for timing control. Secondly, to address the computational inaccuracy caused by the limited conductance range of memristors, we develop an optimized matrix-to-conductance mapping scheme that jointly considers device physical constraints and matrix statistics, achieving over 60% reduction in relative computation error compared with baseline scheme. An associated lightweight circuit enhancement ensures compatibility with practical crossbar architectures, without incurring significant hardware overhead. Thirdly, we establish a memristor programming time model grounded in device-level potentiation and depression dynamics. The analysis yields closed-form expressions for the expected programming time and its upper bound, and is further validated through Monte Carlo simulations, enabling accurate estimation of the system throughput. Simulation results demonstrate that the

proposed circuit achieves a bit error rate comparable to that of 64-bit floating-point precoding, while delivering over 100× improvement in both energy and area efficiency compared with the NVIDIA RTX A2000 graphics processing unit (GPU).

**Index Terms**—Massive MIMO, precoding, ZF and MMSE integration, memristor crossbar, in-memory computing.

## I. INTRODUCTION

MASSIVE multiple-input multiple-output (MIMO) technology has become a key feature of modern cellular communication systems, as exemplified by the 5G cellular networks [1]. Compared with conventional small-scale MIMO systems, massive MIMO offers superior interference mitigation capabilities and higher spectral efficiency. In the 6G mobile networks, massive MIMO is expected to scale up to hundreds or even thousands of antennas to further enhance system performance.

In time-division duplexing (TDD) systems, precoding techniques utilize uplink channel state information (CSI) to preprocess downlink signals, effectively reducing inter-user interference [2]. Common linear precoding algorithms include maximum ratio transmission (MRT), zero-forcing (ZF) precoding, and minimum mean square error (MMSE) precoding. Notably, ZF and MMSE algorithms involve matrix inversion operations with cubic computational complexity [3], leading to significant increases in processing latency and energy consumption in the baseband units of massive MIMO systems.

The emergence of memristor crossbars offers a promising approach to accelerating matrix operations. Memristors can perform matrix-vector multiplications (MVM) with in-memory computing architecture [4], leveraging highly parallel processing to substantially reduce computational complexity, energy consumption, and data transfer latency between memory and processors [5].

Previous studies explored the use of memristor crossbar circuits to implement ZF and regularized ZF algorithms based on ridge regression [6]. Later research designed memristor-accelerated ZF precoding with timing control,

Received 1 August 2025; revised 17 November 2025; accepted 17 December 2025. Date of publication 29 December 2025; date of current version 15 June 2026. This work was supported in part by the National Natural Science Foundation under Grant 62550173 and in part by Beijing Municipal Natural Science Foundation under Grant L242013. This article was recommended by Guest Editor X. Li. (*Corresponding author: Shaoshi Yang.*)

Yu-Xin Zhang is with the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China, also with the Key Laboratory of Mathematics and Information Networks, Ministry of Education, Beijing 100876, China, and also with the Department of IoT Technologies and Applications, China Mobile Research Institute, Beijing 100053, China (e-mail: yuxin.zhang@bupt.edu.cn).

Shaoshi Yang and Yi-Hang Ren are with the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China, and also with the Key Laboratory of Mathematics and Information Networks, Ministry of Education, Beijing 100876, China (e-mail: shaoshi.yang@bupt.edu.cn; renyihang@bupt.edu.cn).

Sheng Chen is with the School of Electronics and Computer Science, University of Southampton, SO17 1BJ Southampton, U.K. (e-mail: sqc@ecs.soton.ac.uk).

Ping Zhang is with the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China, and also with the State Key Laboratory of Networking and Switching Technology, Beijing 100876, China (e-mail: pzhang@bupt.edu.cn).

Digital Object Identifier 10.1109/JETCAS.2025.3648808

2156-3357 © 2025 IEEE. All rights reserved, including rights for text and data mining, and training of artificial intelligence and similar technologies. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

Authorized licensed use limited to: UNIVERSITY OF SOUTHAMPTON. Downloaded on June 18, 2026 at 15:52:51 UTC from IEEE Xplore. Restrictions apply.

where the matrix inversion (INV) and matrix-vector multiplication processes are separated and executed sequentially [7]. Efforts also focused on employing memristor crossbars for high-parallelism maximum likelihood (ML) detection in massive MIMO systems [8]. More recently, hybrid analog-digital computing architectures were introduced to develop memristor-based circuits for successive interference cancellation (SIC) detection in massive MIMO systems [9]. To further enhance robustness against conductance deviations, an amplifier-augmented detector architecture was proposed to decouple the processing of large-scale and small-scale fading matrices, achieving improved bit error rate (BER) with negligible power overhead increase [10]. A follow-up work further introduced refined mathematical modeling and conductance mapping schemes to enhance detection accuracy under memristor deviations [11].

Despite these advancements, these prior works have not provided a detailed assessment of how memristor non-idealities affect precoding performance. For instance, the work [7] focused on overall BER trends under noise, but did not explicitly isolate the impact of conductance quantization or programming uncertainty. Similarly, the study [6] emphasized errors under varying amplifier gain and memory precision, while only offering very limited analysis on how circuit performance is affected by quantization constraints or practical conductance ranges in low-bit precision regimes.

Moreover, few studies have systematically assessed how memristor non-idealities, such as limited switching ratio (the ratio of high resistance state to low resistance state) and programming variability, affect the performance of precoding circuits. Notably, the work [12] proposed a resistive random access memory (RRAM) based baseband processor for MIMO orthogonal frequency-division multiplexing (OFDM) systems and introduced a model to estimate the row-wise programming time of memristor crossbars. Their analysis reveals that the latency scales sublinearly with the number of transmit antennas. While their work marks a valuable step toward programming time modeling, it adopts simplified conductance-programming step model and focuses mainly on asymptotic upper-bound estimation, without considering the impact of realistic programming behavior on the latency.

To bridge these gaps, we propose a memristor crossbar circuit that supports both ZF and MMSE precoding. The proposed circuit reduces energy consumption by minimizing the crossbar size and peripheral circuitry, and it enables fast one-step computation without requiring timing control. To complement the circuit design, we further develop a probabilistic model, supported by both theoretical analysis and simulation results, to more accurately evaluate the memristor programming time under realistic device behaviors, such as potentiation and depression dynamics.

These contributions bridge the modeling gap between physical device constraints and system-level performance, enabling more accurate design trade-offs for future memristor-based massive MIMO baseband processors. The main contributions of this paper are summarized as follows.

- We propose a one-step memristor crossbar circuit that supports both ZF and MMSE precoding, featuring a

reduced array size and simplified peripheral circuitry. By eliminating the need for timing control, the circuit further accelerates the processing speed through efficient one-step computation.

- To mitigate the degradation of computational accuracy caused by the limited switching ratio of memristors, we develop an optimized matrix-to-conductance mapping scheme that jointly considers the physical constraints of memristors and the statistical properties of the precoding matrix. Furthermore, the high conductance values assigned to diagonal elements, as dictated by the proposed mapping scheme, motivate a resistor-parallel circuit architecture selectively applied to diagonal memristor cells via metal-oxide-semiconductor field-effect transistor (MOSFET) regulated switching, thereby minimizing hardware overhead. The proposed mapping scheme achieves over 60% reduction in relative computation error compared with the baseline linear mapping scheme, demonstrating superior robustness under various memristor physical constraints.
- We propose a probabilistic model for estimating the memristor programming time, which accounts for the potentiation and depression characteristics of memristors. Based on this model, we derive closed-form expressions for the statistical expectation of programming time and its upper bound, thereby providing theoretical support for the time estimation. The consistency between the theoretical predictions and simulation results validates the accuracy of the proposed model. This estimation scheme offers a reliable approach for the quantitative analysis of circuit throughput, energy efficiency and area efficiency. Based on the theoretical analysis and simulation results, the throughput is comparable to the advanced commercial processors and both the energy and area efficiency of our design is 100 times higher than these processors.

## II. PRELIMINARIES

### A. System Model

The transmitter considered in this paper is located at the base station (BS) of a massive MIMO system operating in the downlink scenario. The channel matrix  $\tilde{\mathbf{H}} \in \mathbb{C}^{N_{\text{rx}} \times N_{\text{tx}}}$  is used to characterize the channel gains between the BS and the user equipment (UE), where  $N_{\text{tx}}$  denotes the number of transmit antennas at the BS and  $N_{\text{rx}}$  denotes the number of receive antennas at the UE. The tilde symbol, e.g.,  $\tilde{\mathbf{H}}$ , is used to indicate complex-valued variables throughout the paper. Each element  $\tilde{h}_{n_{\text{rx}}, n_{\text{tx}}}$  of the channel matrix  $\tilde{\mathbf{H}}$  represents the subchannel gain between the  $n_{\text{tx}}$ -th transmit antenna at the BS and the  $n_{\text{rx}}$ -th receive antenna at the UE.

Under ideal conditions, where the distance between transmit and receive antennas is sufficiently large and the channel environment is highly random, the subchannel gains in the MIMO system can be assumed to be uncorrelated. Such channel conditions can be modeled by an independent and identically distributed Rayleigh fading channel, which is  $\tilde{h}_{n_{\text{rx}}, n_{\text{tx}}} \sim \mathcal{CN}(0, 1)$ .

The received signals at the  $N_{\text{rx}}$  antennas of the UE are given by:

$$\tilde{\mathbf{y}} = \sqrt{\rho_T} \tilde{\mathbf{H}} \tilde{\mathbf{x}} + \tilde{\mathbf{n}}, \quad (1)$$

where  $\sqrt{\rho_T}$  denotes the transmit power factor,  $\tilde{\mathbf{x}}$  represents the transmitted complex signal vector from the  $N_{\text{tx}}$  antennas, and  $\tilde{\mathbf{n}} \sim \mathcal{CN}(0, \sigma_n^2 \mathbf{I}_{N_{\text{rx}}})$  denotes the additive white Gaussian noise with  $\mathbf{I}_{N_{\text{rx}}}$  being the  $N_{\text{rx}} \times N_{\text{rx}}$  identity matrix.

### B. Precoding Algorithms

The linear precoding operation is given by  $\tilde{\mathbf{x}} = \tilde{\mathbf{W}} \tilde{\mathbf{s}}$ , where  $\tilde{\mathbf{W}}$  denotes the precoding matrix and  $\tilde{\mathbf{s}}$  denotes the transmitted symbol vector.

ZF precoding aims to direct the transmitted signal energy toward the intended user while minimizing inter-antenna interference by forcing the signals in other directions to near zero [13]. The ZF precoding matrix is given by:

$$\tilde{\mathbf{W}} = \gamma \tilde{\mathbf{H}}^H (\tilde{\mathbf{H}} \tilde{\mathbf{H}}^H)^{-1}, \quad (2)$$

where  $\gamma$  is a scalar factor ensuring the transmit power constraint.

MMSE precoding aims to minimize the mean square error between the transmitted and received signals by considering both inter-antenna interference and noise [14]. The MMSE precoding matrix is given by:

$$\tilde{\mathbf{W}} = \gamma \tilde{\mathbf{H}}^H (\tilde{\mathbf{H}} \tilde{\mathbf{H}}^H + \lambda \mathbf{I}_{N_{\text{tx}}})^{-1}, \quad (3)$$

where  $\lambda = \frac{\sigma_s^2}{\sigma_n^2}$  is the ratio of signal power to noise power, which can be regarded as a regularization parameter.

Considering the physical constraints that conductance, current, and voltage are real-valued in memristor crossbar circuits, complex-valued vectors/matrices must be converted into real-valued forms, in order to perform computation within memristor crossbar circuits. For vector  $\tilde{\mathbf{v}}$  and matrix  $\tilde{\mathbf{M}}$ , their corresponding real-valued forms are given by

$$\mathbf{\Omega}_v \triangleq [\Re(\tilde{\mathbf{v}})^T \quad \Im(\tilde{\mathbf{v}})^T]^T, \quad \mathbf{\Omega}_M \triangleq \begin{bmatrix} \Re(\tilde{\mathbf{M}}) & -\Im(\tilde{\mathbf{M}}) \\ \Im(\tilde{\mathbf{M}}) & \Re(\tilde{\mathbf{M}}) \end{bmatrix}, \quad (4)$$

where  $\Re(\cdot)$  and  $\Im(\cdot)$  denote the real and imaginary parts, respectively. We also define  $\mathbf{x}^+ \triangleq \max\{x, 0\}$  and  $\mathbf{x}^- \triangleq -\min\{x, 0\}$ .

### C. Basic Characteristics of Memristors

The minimum and maximum available conductance values of the memristor are denoted by  $G_{\min}$  and  $G_{\max}$ . By applying appropriate electrical excitation, the conductance can be adjusted to a target value  $G_{\text{target}} \in [G_{\min}, G_{\max}]$ . The actual conductance is denoted as  $G_{\text{final}}$ . Experimental results have shown that the conductance variation follows a zero-mean Gaussian distribution with a variance that remains constant regardless of the conductance value [15]. To simplify the analysis, the memristor conductance variation is modeled as  $G_{\text{final}} = \delta G_{\text{target}}$ , where  $\delta$  denotes the conductance variation factor and  $\delta \sim \mathcal{N}(1, \sigma_{\text{norm}}^2)$ . The normalized standard deviation is given by  $\sigma_{\text{norm}} = \epsilon_p / G_{\text{target}}$ , where  $\epsilon_p$  represents the small residual error after programming process, and it has the same unit as  $G_{\text{target}}$ .

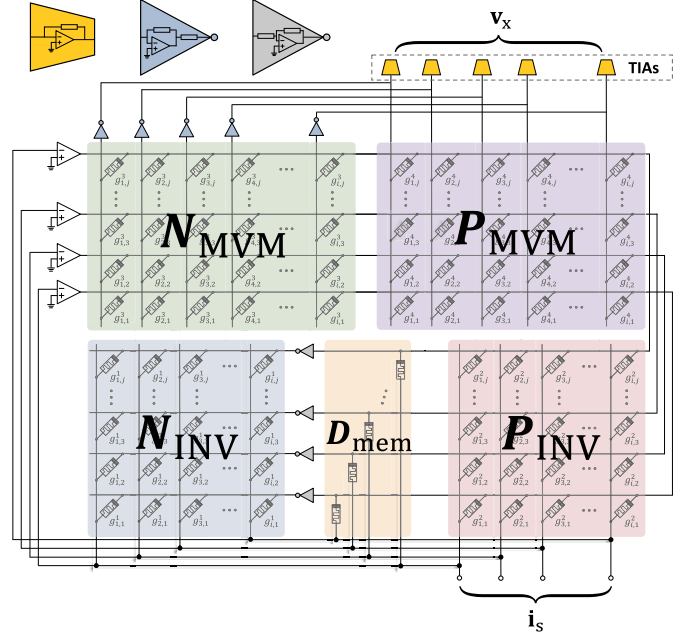


Fig. 1. Proposed memristor crossbar-based MIMO ZF & MMSE precoder circuit.

The conductance resolution is represented by bit precision, which specifies the number of resolvable conductance levels, and it is influenced by the device variation during programming process. In practical implementations, considering the limitation of programming time, the programming process relies on a predefined lookup table, and the target conductance  $G_{\text{target}}$  is restricted to a discrete set of values rather than being continuously tunable. Typically, the conductance range is divided into  $L = 2^{\epsilon_{\text{bp}}}$  levels [16], and  $\epsilon_{\text{bp}}$  denotes the bit precision of the memristor.

### D. Basic Circuit Theory for Memristor Crossbar

1) *Matrix-Vector Multiplication*: The MVM circuit executes the operation  $\mathbf{M} \cdot \mathbf{s} = \mathbf{n}$ . Its core component is a memristor crossbar consisting of  $i$  word lines (rows) and  $j$  bit lines (columns). The conductance at each crosspoint forms the conductance matrix  $\mathbf{G} \in \mathbb{R}^{i \times j}$ , which is linearly mapped from  $\mathbf{M} \in \mathbb{R}^{i \times j}$ . According to Kirchhoff's current law, the relationship between the output current vector  $\mathbf{i}$ , the conductance matrix  $\mathbf{G}$ , and the input voltage vector  $\mathbf{v}$  is given by:

$$\mathbf{i} = \mathbf{G} \cdot \mathbf{v}. \quad (5)$$

2) *Matrix Inversion*: The matrix inversion circuit shares a similar structure with the MVM circuit but operates on square matrices. By applying Kirchhoff's laws, the relationship is given as:

$$-\mathbf{G}^{-1} \cdot \mathbf{i} = \mathbf{v}. \quad (6)$$

This circuit can also be regarded as a linear equation solver, capable of computing  $\mathbf{M}^{-1} \cdot \mathbf{s} = \mathbf{n}$  in a single step, which is particularly valuable for channel inversion in linear precoding.

### III. PROPOSED MEMRISTOR CROSSBAR-BASED MIMO ZF AND MMSE PRECODER CIRCUIT

The circuit consists of two main components: the INV component (including  $\mathbf{N}_{\text{INV}}$ ,  $\mathbf{D}_{\text{mem}}$  and  $\mathbf{P}_{\text{INV}}$ ) and the MVM component (including  $\mathbf{N}_{\text{MVM}}$  and  $\mathbf{P}_{\text{MVM}}$ ), as illustrated in Fig. 1. In the INV component, the diagonal matrix  $\mathbf{D}_{\text{mem}}$  can be regarded as the sum of a balancing matrix and a regularization matrix. The balancing matrix addresses the conductance imbalance caused by the diagonally dominant matrix  $\widetilde{\mathbf{H}}\widetilde{\mathbf{H}}^H$ , while the regularization matrix maps the regularization term ( $\lambda\mathbf{I}$ ) in the MMSE algorithm. In the MVM component, the operational transconductance amplifiers (OTAs) convert the output voltage  $\mathbf{v}_x$  to current feedback. By applying Ohm's law and Kirchhoff's laws to the circuit shown in Fig. 1, the following equation is obtained:

$$(\mathbf{P}_{\text{INV}} - \mathbf{N}_{\text{INV}} + \mathbf{D}_{\text{mem}})^{-1}\mathbf{i}_s = (\mathbf{P}_{\text{MVM}} - \mathbf{N}_{\text{MVM}})^{-1}\mathbf{v}_x. \quad (7)$$

Denote  $\widetilde{\mathbf{Z}} = \widetilde{\mathbf{H}}\widetilde{\mathbf{H}}^H$  and construct the matrix  $\mathbf{\Omega}_M = \frac{2\mathbf{\Omega}_Z}{N_{\text{tx}}} - 2\mathbf{I}_{2N_{\text{rx}}}$ . Then the proposed circuit mapping scheme is given as in (8), shown at the bottom of the page, where the factor  $\alpha$  is employed to establish the correspondence between matrix elements and memristor conductance values.

For  $\lambda = 0$ , the proposed circuit performs the ZF precoding algorithm, while for  $\lambda \neq 0$ , the proposed circuit performs the MMSE precoding algorithm.

### IV. OPTIMIZED MAPPING SCHEME FOR MIMO PRECODER CIRCUIT

It is widely acknowledged that memristors cannot directly represent negative values, and the adopted positive-negative separation scheme results in approximately 50% of the target conductance values being zero. Consequently, limited by the memristor switching ratio, different mapping scales can significantly affect computational errors. In fact, through a series of appropriate linear scaling operations, the proposed circuit retains the capability to accurately perform the intended computation. Based on the circuit structure, a more general mapping expression can be derived as in (9), shown at the bottom of the page, where the factor  $\kappa$  has the similar role to  $\alpha$ ,  $r_{\text{ds}} = \frac{N_{\text{rx}}}{N_{\text{d}}}$  represents the matrix scaling factor, and  $N_{\text{d}}$  denotes the scalar that corresponds to the identical diagonal entries of the balancing matrix.

We simulate the computation errors under different mapping ratios and maximum conductance to verify the effectiveness of the proposed mapping scheme. In the simulation,  $\alpha$  is set to

$100\mu\text{S}$ ,  $\kappa$  is set to  $\frac{r_{\text{ds}}G_{\text{max}}}{2\sqrt{2}}$ ,  $G_{\text{max}}$  is expressed in unit of  $100\mu\text{S}$ ,  $\mathbf{i}_s$  is expressed in milliamperes, and  $\mathbf{v}_x$  is expressed in millivolts. The simulation results are shown in Fig. 2(a), which indicate that the relative computation errors of the proposed circuit decreases as  $N_{\text{d}}$  increases within a certain range. However, when  $N_{\text{d}}$  exceeds this range, the relative error increases rapidly. This is because limited by the conductance range of the memristor, an excessively large  $N_{\text{d}}$  significantly raises the probability of conductance exceeding the conductance range, leading to computational inaccuracies. As  $G_{\text{max}}$  increases, the feasible range of  $N_{\text{d}}$  also expands, which facilitates a further reduction in computation error. Therefore, when selecting the mapping ratio, a trade-off should be considered based on the characteristics of the memristor device. The simulation results demonstrate that the optimized mapping scheme reduces the relative computation error by more than 60% compared to the baseline approaches ( $N_{\text{d}} = 2$ ) in [7], demonstrating superior robustness across various memristor physical constraints.

To ensure that the probability of conductance exceeding the limitation range of memristor conductance remains below 0.3%, the theoretical optimal value of  $N_{\text{d}}$  can be determined through calculation, and the optimal value  $N_{\text{d}}^*$  is given by

$$N_{\text{d}}^* = \xi \cdot \frac{\sqrt{2N_{\text{tx}}}}{3} \cdot G_{\text{max}}, \quad (10)$$

where  $\xi$  is an experimentally determined coefficient to compensate for the random disturbances, such as residual error  $\epsilon_p$ , that affect circuit performance. In the simulation,  $\xi$  is set to 0.8.

To further validate the robustness of the proposed mapping scheme under spatially correlated MIMO channels, additional simulations are conducted based on the Kronecker channel model, where the channel matrix is given by  $\widetilde{\mathbf{H}} = \mathbf{R}_t^{1/2}\widetilde{\mathbf{W}}\mathbf{R}_r^{1/2}$ . In this model,  $\widetilde{\mathbf{W}}$  denotes a Rayleigh fading matrix with i.i.d. complex Gaussian entries, and both the transmit and receive correlation matrices  $\mathbf{R}_t$  and  $\mathbf{R}_r$  follow a Toeplitz exponential model. The  $(i, j)$ -th entry of  $\mathbf{R}$  is given by  $[\mathbf{R}]_{i,j} = \rho^{|i-j|}$ . The transmit and receive correlation coefficients are assumed to be identical ( $\rho_t = \rho_r = \rho$ ).

The corresponding results are shown in Figs. 2(b)–(d), which illustrate the variation of relative computation error with respect to  $N_{\text{d}}$  for different correlation coefficients  $\rho$ , under maximum conductance of  $200\mu\text{S}$ ,  $300\mu\text{S}$ , and  $400\mu\text{S}$ , respectively. Similar to the Rayleigh case, the relative error first decreases and then increases as  $N_{\text{d}}$  grows. In addition, higher spatial correlation not only results in slightly larger

$$\begin{aligned} \mathbf{i}_s &= -\mathbf{\Omega}_s, \quad \mathbf{v}_x = \mathbf{\Omega}_x, \quad \mathbf{D}_{\text{mem}} = \alpha \left( 2 + \frac{2\lambda}{N_{\text{tx}}} \right) \mathbf{I}_{2N_{\text{rx}}}, \quad \mathbf{P}_{\text{INV}} = \alpha \mathbf{\Omega}_M^+, \\ \mathbf{N}_{\text{INV}} &= \alpha \mathbf{\Omega}_M^-, \quad \mathbf{P}_{\text{MVM}} = \alpha \frac{2\mathbf{\Omega}_{\text{HH}}^+}{N_{\text{tx}}}, \quad \mathbf{N}_{\text{MVM}} = \alpha \frac{2\mathbf{\Omega}_{\text{HH}}^-}{N_{\text{tx}}}, \end{aligned} \quad (8)$$

$$\begin{aligned} \mathbf{i}_s &= -\frac{\mathbf{\Omega}_s}{\kappa}, \quad \mathbf{v}_x = \mathbf{\Omega}_x, \quad \mathbf{D}_{\text{mem}} = \alpha \left( N_{\text{d}} + \frac{\lambda}{r_{\text{ds}}} \right) \mathbf{I}_{2N_{\text{rx}}}, \quad \mathbf{P}_{\text{INV}} = \alpha \left( \frac{\mathbf{\Omega}_Z}{r_{\text{ds}}} - N_{\text{d}} \mathbf{I}_{2N_{\text{rx}}} \right)^+, \\ \mathbf{N}_{\text{INV}} &= \alpha \left( \frac{\mathbf{\Omega}_Z}{r_{\text{ds}}} - N_{\text{d}} \mathbf{I}_{2N_{\text{rx}}} \right)^-, \quad \mathbf{P}_{\text{MVM}} = \frac{\kappa \mathbf{\Omega}_{\text{HH}}^+}{r_{\text{ds}}}, \quad \mathbf{N}_{\text{MVM}} = \frac{\kappa \mathbf{\Omega}_{\text{HH}}^-}{r_{\text{ds}}}, \end{aligned} \quad (9)$$

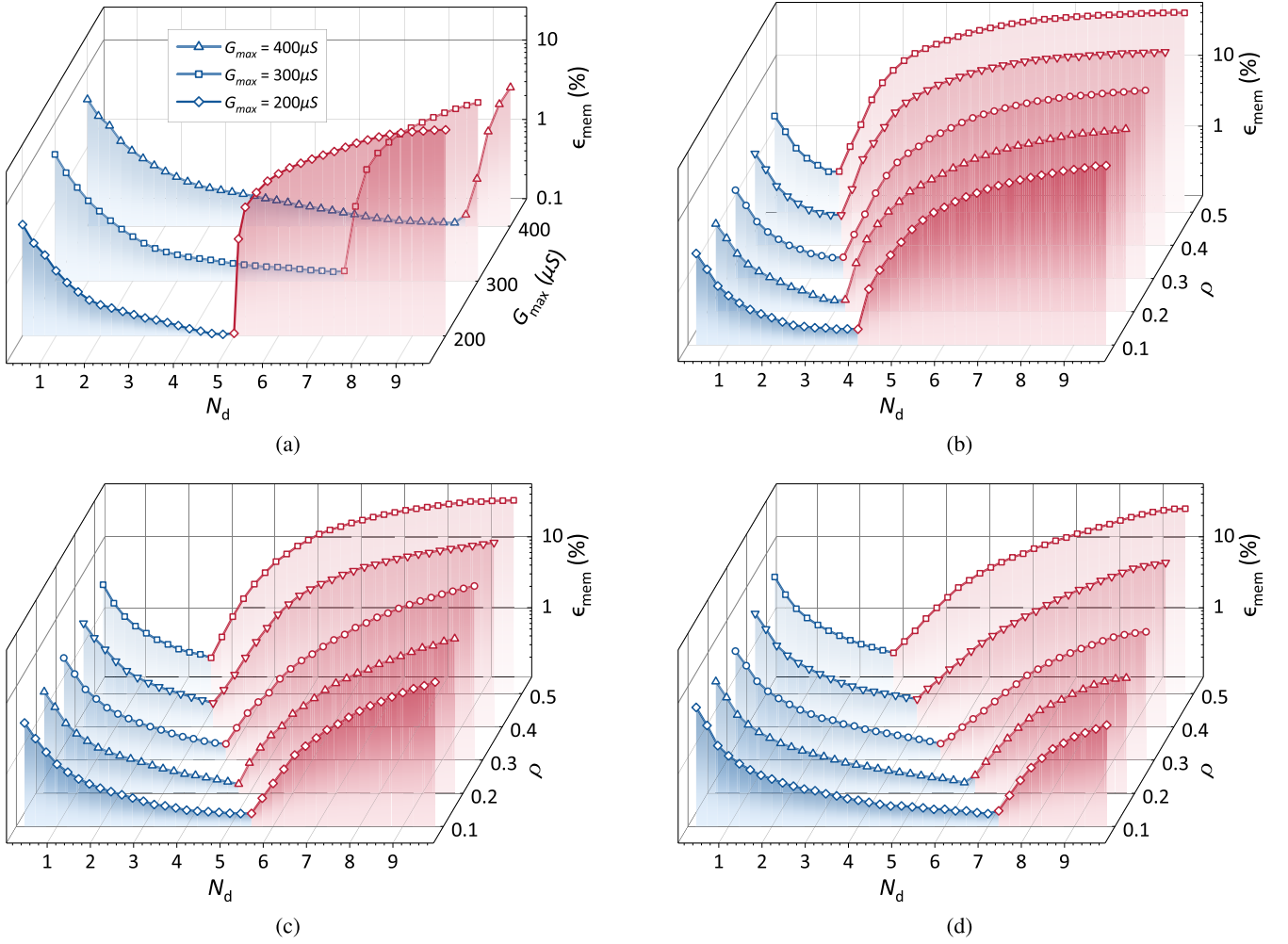


Fig. 2. Circuit computation errors under different mapping ratios and maximum conductance (with same minimum conductance). (a) Rayleigh fading channel, (b) Kronecker channel model ( $G_{\max} = 200\mu\text{S}$ ), (c) Kronecker channel model ( $G_{\max} = 300\mu\text{S}$ ), (d) Kronecker channel model ( $G_{\max} = 400\mu\text{S}$ ).

computation errors due to the uneven distribution of equivalent channel gains, but also narrows the feasible range of  $N_d$ . Consequently, although increasing  $G_{\max}$  still enlarges the valid  $N_d$  range and helps further reduce the overall error, the benefit gradually diminishes under stronger correlation. These results demonstrate that the proposed mapping scheme remains reasonably effective under correlated channel conditions, particularly in relatively rich-scattering channels.

Similar to the Rayleigh channel case, for the Kronecker channel model, the optimal value  $N_d^*$  is expressed as

$$N_d^* = \xi \cdot \frac{G_{\max} N_t}{\eta\rho + 3\sqrt{\frac{\zeta}{2}}(1 + \rho)}, \quad (11)$$

where  $\eta \triangleq \text{tr}(R_t)$  and  $\zeta \triangleq \text{tr}(R_t^2)$  are determined by the transmit correlation matrix  $R_t$ . Compared with the Rayleigh channel case, the presence of spatial correlation increases the denominator term in (11), thereby reducing the feasible range of  $N_d$  and resulting in a smaller theoretical optimum  $N_d^*$ . This analytical result is consistent with the simulation observations in Figs. 2(b)–(d), confirming that higher correlation levels constrain the allowable mapping ratio and thereby affect the achievable computation accuracy.

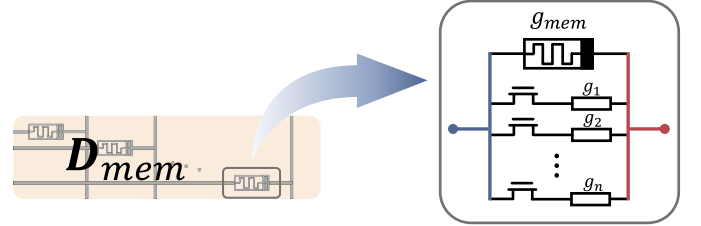


Fig. 3. Circuit structure of one diagonal memristor cell.

Unlike the target conductance of the memristors in the MVM and INV components shown in Fig. 1, the diagonal matrix  $\mathbf{D}_{\text{mem}}$  corresponds to a relatively higher set of memristor conductance values. To accommodate this feature while minimizing hardware overhead, a lightweight circuit enhancement is introduced: upon considering the limited number of diagonal elements in  $\mathbf{D}_{\text{mem}}$ , each memristor in the diagonal matrix is connected in parallel with multiple resistors of fixed conductance  $G_{\max}$ , while MOSFET switches are utilized to regulate their connection to the circuit and flexibly realize the desired conductance values. The corresponding circuit structure is illustrated in Fig. 3.

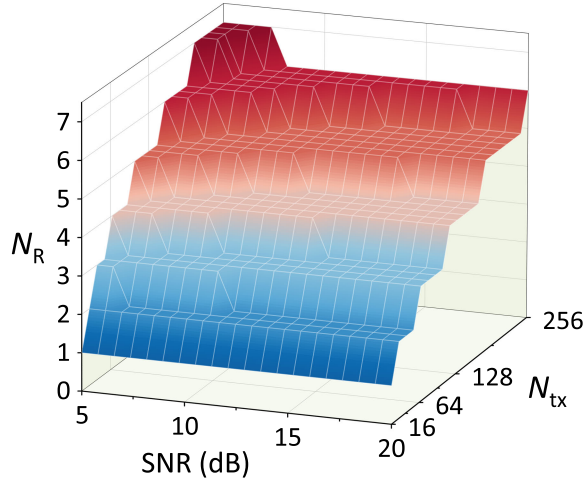


Fig. 4. Number of constant resistors in each memristor cell under Rayleigh channel.

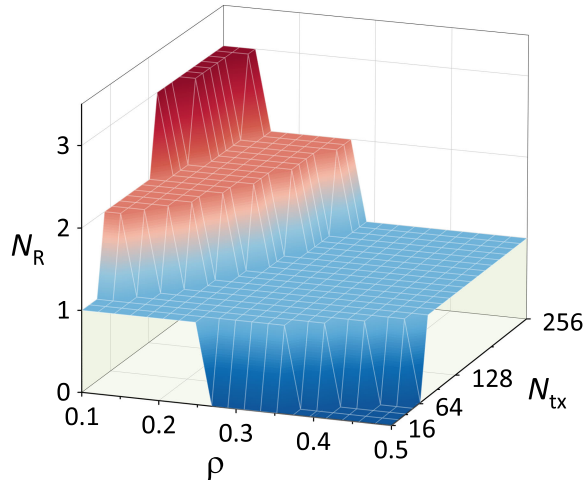


Fig. 5. Number of constant resistors in each memristor cell under Kronecker channel model, with SNR fixed at 10 dB.

Under the assumption of Rayleigh channel, the number of constant resistors  $N_R$  contained in each memristor cell can be determined according to

$$N_R = \left\lceil \xi \cdot \left( \frac{\lambda}{N_{tx}} + 1 \right) \frac{\sqrt{2N_{tx}}}{3} \right\rceil, \quad (12)$$

where  $\lceil \cdot \rceil$  denotes the ceiling function. The value of  $N_R$  primarily depends on the memristor crossbar scale and is relatively less sensitive to the SNR. Fig. 4 illustrates the number of constant resistors required per memristor cell in the proposed circuit when executing the MMSE precoding algorithm under different crossbar scales and SNR conditions.

When spatial correlation is introduced in the Kronecker channel model, the estimation of  $N_R$  is modified as

$$N_R = \left\lceil \xi \cdot \left( \frac{\lambda}{N_{tx}} + 1 \right) \cdot \frac{N_t}{\eta\rho + 3 \sqrt{\frac{\xi}{2}(1+\rho)}} \right\rceil. \quad (13)$$

Fig. 5 further illustrates this relationship by fixing the SNR at 10 dB and plotting the variation of  $N_R$  with respect to different crossbar scales and correlation coefficients.

## V. A MONTE CARLO APPROACH FOR MEMRISTOR PROGRAMMING TIME ESTIMATION

### A. Closed-Form Expression of Programming Time

Considering the proposed circuit in the downlink scenario, programming time constitutes a significant portion of the overall computational latency. However, to the best of our knowledge, there has been few prior works that provide a detailed analysis of the programming time. In this subsection, we derive a closed-form expression for the memristor programming time and validate it through Monte Carlo simulations.

Throughout the derivation,  $w$  denotes the normalized ratio of programming steps,  $T_s$  represent the duration of a single programming pulse, and  $S_{total}$  denote the total number of steps required to switch a memristor from its lowest to highest resistance state. The memristor crossbar commonly implements row-wise parallel programming [8]. Since the impact of residual errors on the programming time is relatively minor, they are excluded in the following simplified model.

Building on the preceding discussion, we now derive a closed-form expression for the expected number of programming steps. We take the matrix  $\mathbf{P}_{INV}$  as an example for derivation. The off-diagonal elements of this matrix follow a Gaussian distribution, while the diagonal elements follow a centralized Gamma distribution. Based on this distinction, the modeling of diagonal and off-diagonal elements is carried out separately as follows.

**Off-diagonal:** Each off-diagonal element is modeled as a zero-mean Gaussian random variable:

$$P \sim \mathcal{N}(0, \sigma^2), \quad \sigma = \alpha \frac{N_d}{\sqrt{2N_{tx}}}. \quad (14)$$

**Diagonal:** Each diagonal element is modeled as a centralized Gamma random variable with a scaling factor:

$$P \sim \Gamma(k = N_t, \theta) = \alpha N_d, \quad \theta = \alpha \frac{2N_d}{N_{tx}}. \quad (15)$$

All conductance values are represented by the same variable  $P$  for the sake of illustration, and the unit of  $P$  is  $\mu\text{S}$ . The conductance variable  $P$  is quantized over a predefined range:

$$\begin{cases} G_k = G_{\min} + k\Delta G, & k = 0, 1, \dots, L-1, \\ \Delta G = \frac{G_{\max} - G_{\min}}{L}. \end{cases} \quad (16)$$

The quantization function is defined as:

$$Q(P) = \begin{cases} G_0, & P \leq G_0, \\ G_k, & G_k < P \leq G_{k+1}, \quad 1 \leq k \leq L-1, \\ G_{L-1}, & P > G_{L-1}. \end{cases} \quad (17)$$

To more accurately estimate the programming time of memristors, we adopt the programming model proposed in [17]. Let  $\alpha_p$  and  $\alpha_d$  represent the potentiation and depression

coefficients respectively, which characterize the nonlinearity of the memristor programming process. As shown in (18) and (19), the functions  $G_{\text{pot}}(w)$  and  $G_{\text{dep}}(w)$  characterize the evolution of memristor conductance with respect to the normalized number of programming steps under the potentiation and depression operations, respectively.

$$G_{\text{pot}}(w) = \left( (G_{\text{max}}^{\alpha_p} - G_{\text{min}}^{\alpha_p}) w + G_{\text{min}}^{\alpha_p} \right)^{\frac{1}{\alpha_p}}, \quad w \in [0, 1], \quad (18)$$

$$G_{\text{dep}}(w) = \left( (G_{\text{max}}^{\alpha_d} - G_{\text{min}}^{\alpha_d}) w + G_{\text{min}}^{\alpha_d} \right)^{\frac{1}{\alpha_d}}, \quad w \in [0, 1]. \quad (19)$$

From these two equations, we obtain the inverse function of (18) and (19):

$$w(G) = \begin{cases} \frac{G^{\alpha_p} - G_{\text{min}}^{\alpha_p}}{G_{\text{max}}^{\alpha_p} - G_{\text{min}}^{\alpha_p}}, & G_{\text{tar}} > G_{\text{cur}}, \\ 1 - \frac{G^{\alpha_d} - G_{\text{min}}^{\alpha_d}}{G_{\text{max}}^{\alpha_d} - G_{\text{min}}^{\alpha_d}}, & G_{\text{tar}} \leq G_{\text{cur}}, \end{cases} \quad (20)$$

where  $G_{\text{cur}}$  denotes the current conductance, and  $G_{\text{tar}}$  denotes the target conductance. From  $G_{\text{cur}}$  to  $G_{\text{tar}}$ , the required number of programming steps is given by

$$S_{\text{prog}} = S_{\text{total}} \cdot |w(G_{\text{tar}}) - w(G_{\text{cur}})|. \quad (21)$$

For each quantized level  $G_k$ , the corresponding probability mass function,  $p_k = \Pr(G_k \leq P < G_{k+1})$ , is detailed below.

**Off-diagonal:**

$$p_k = \begin{cases} \Phi\left(\frac{G_1}{\sigma}\right), & k = 0, \\ \Phi\left(\frac{G_{k+1}}{\sigma}\right) - \Phi\left(\frac{G_k}{\sigma}\right), & 1 \leq k \leq L-2, \\ 1 - \Phi\left(\frac{G_{L-1}}{\sigma}\right), & k = L-1. \end{cases} \quad (22)$$

**Diagonal:**

$$p_k = \begin{cases} F(G_1 + \alpha N_d), & k = 0, \\ F(G_{k+1} + \alpha N_d) - F(G_k + \alpha N_d), & 1 \leq k \leq L-2, \\ 1 - F(G_{L-1} + \alpha N_d), & k = L-1. \end{cases} \quad (23)$$

where  $\Phi(\cdot)$  and  $F(\cdot)$  are the cumulative distribution functions of the Gaussian and Gamma distributions, respectively. The probabilities are normalized

$$p_k \leftarrow \frac{p_k}{\sum_{i=0}^{L-1} p_i}, \quad (24)$$

to ensure  $\sum_{k=0}^{L-1} p_k = 1$ .

Assuming independently sampled conductance  $G_m$  and  $G_k$ , the expected programming step is given by:

$$\mathbb{E}[S_{\text{prog}}] = S_{\text{total}} \sum_{k=0}^{L-1} \sum_{m=0}^{L-1} p_k p_m |w(G_m) - w(G_k)|. \quad (25)$$

As  $L \rightarrow \infty$ , the expected number of programming steps can be expressed as

$$\mathbb{E}[S_{\text{prog}}] = S_{\text{total}} \int \int |w(y) - w(x)| f_P(x) f_P(y) dx dy, \quad (26)$$

where the probability density function (PDF)  $f_P(x)$  is defined over the interval  $[G_{\text{min}}, G_{\text{max}}]$ . The PDFs are given as follows.

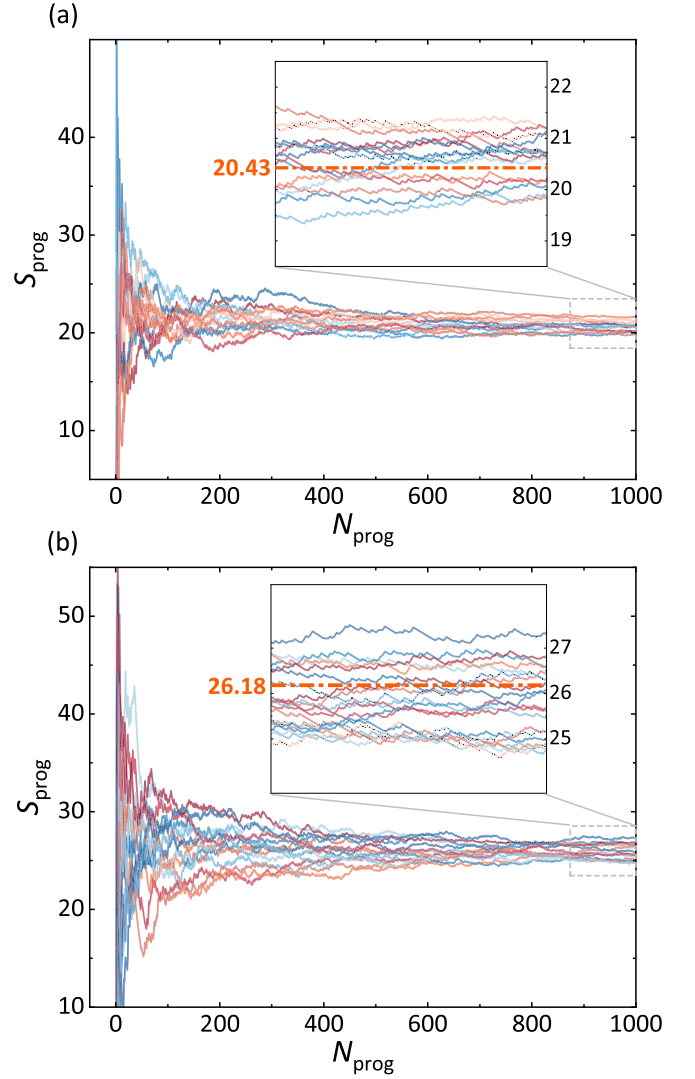


Fig. 6. Convergence of the average number of programming steps obtained via Monte Carlo simulations: (a)  $\mathbf{P}_{\text{INV}}$  matrix, and (b)  $\mathbf{P}_{\text{MVM}}$  matrix. Each curve represents a different initial conductance state. In each subfigure, the theoretical value is shown in orange for comparison.

**Off-diagonal:**

$$f_P(z) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{z^2}{2\sigma^2}\right), \quad (27)$$

**Diagonal:**

$$f_P(z) = \frac{1}{\alpha\Gamma(k)\theta^k} \left(\frac{z + \alpha N_d}{\alpha}\right)^{k-1} \exp\left(-\frac{z + \alpha N_d}{\alpha\theta}\right). \quad (28)$$

The programming time  $T_{\text{prog}}$  of a single memristor is given by:

$$T_{\text{prog}} = T_s \cdot \bar{S}_{\text{prog}} = T_s \cdot \mathbb{E}[S_{\text{prog}}]. \quad (29)$$

### B. Validation by Monte Carlo Simulation

To validate the above expression, we conducted a Monte Carlo simulation, as illustrated in Fig. 6, where the theoretical values are highlighted in orange. The simulation parameters are set as  $T_s = 1$  ns,  $S_{\text{total}} = 100$ ,  $N_{\text{ix}} = 32$  and  $N_d = N_d^*$ .

For each curve, each point represents the average number of programming steps over the first  $N_{\text{prog}}$  independent experiments. As the number of independent programming experiments increases,  $\bar{S}_{\text{prog}}$  gradually stabilizes. When repeating the experiments with different initial conductance states,  $\bar{S}_{\text{prog}}$  consistently converges with less than 5% deviation from the theoretical value. The Monte Carlo results thus validate that our closed-form solution is accurate.

### C. Crossbar Programming Time

Under the row-wise parallel programming scheme, all memristors within a row are updated simultaneously, while memristor crossbar is sequentially programmed from one row to the next. Hence, the total programming time for the crossbar equals to the number of rows times the programming time per row.

For the  $2N_{\text{rx}} \times 2N_{\text{rx}}$  INV crossbar, each row contains one diagonal memristor (programmed with a latency of  $T_{\text{INV-d}}$ ) and multiple off-diagonal memristors (each programmed in  $T_{\text{INV-od}}$ ). Since all memristors within a row are programmed simultaneously, the row programming time is determined by the largest of  $T_{\text{INV-d}}$  and  $T_{\text{INV-od}}$ . Consequently, the overall programming time of the INV crossbar can be expressed as:

$$T_{\text{INV}} = 2N_{\text{rx}} \cdot \max\{T_{\text{INV-d}}, T_{\text{INV-od}}\}. \quad (30)$$

For the  $2N_{\text{rx}} \times 2N_{\text{rx}}$  MVM crossbar, all elements follow the same distribution, and thus the overall programming time of the MVM crossbar is given by

$$T_{\text{MVM}} = 2N_{\text{rx}} \cdot T_{\text{MVM-row}}, \quad (31)$$

where  $T_{\text{MVM-row}}$  denotes the programming time for each row of MVM crossbar.

The overall programming latency is determined by the slowest crossbar, namely,

$$\begin{aligned} T_{\text{prog}} &= \max\{T_{\text{INV}}, T_{\text{MVM}}\} \\ &= 2N_{\text{rx}} \cdot \max\{T_{\text{INV-d}}, T_{\text{INV-od}}, T_{\text{MVM-row}}\}. \end{aligned} \quad (32)$$

### D. Programming Time of Memristor Crossbar in Massive MIMO Systems

In massive MIMO systems, the impact of memristor variability on programming steps is more significant due to the large number of memristors updating in parallel. While Section V-A focuses on the programming time of a single memristor cell, this alone does not reflect the actual programming latency. Specifically, since all memristors in a row are programmed simultaneously, the row programming latency is determined by the slowest device, i.e., the one requiring the maximum number of steps. This effect becomes more significant as the crossbar size increases. To demonstrate the time complexity, we employ extreme value theory to model the statistical distribution of the maximum programming steps.

We take the programming time of matrix  $\mathbf{P}_{\text{INV}}$  as an example for derivation. As derived in (25), the mean  $\mu_S$  and variance  $\sigma_S^2$  of  $S_{\text{prog}}$  can be obtained through:

$$\mu_S = \mathbb{E}[S_{\text{prog}}], \quad \sigma_S^2 = \mathbb{E}[S_{\text{prog}}^2] - (\mathbb{E}[S_{\text{prog}}])^2. \quad (33)$$

Since  $S_{\text{prog}}$  is bounded between 0 and  $S_{\text{total}}$ , it satisfies a sub-Gaussian tail bound, expressed as

$$\Pr\{S_{\text{prog}} - \mu_S > x\} \leq \exp\left(-\frac{x^2}{2\sigma_S^2}\right), \quad x \geq 0. \quad (34)$$

Assume that there are  $M_{\text{off}} = 2N_t - 1$  off-diagonal memristors in each row. Let  $S_{\text{max}}$  denote the maximum programming steps among these memristors, which is expressed as:

$$S_{\text{max}} = \max_{1 \leq j \leq M_{\text{off}}} S_j, \quad (35)$$

where  $S_j$  are i.i.d. copies of  $S_{\text{prog}}$ . Applying the sub-Gaussian tail inequality (34), the survival function of  $S_{\text{max}}$  is upper bounded by

$$\Pr\{S_{\text{max}} > \mu_S + x\} \leq M \exp\left(-\frac{x^2}{2\sigma_S^2}\right). \quad (36)$$

The expected maximum programming steps can be expressed as

$$\mathbb{E}[S_{\text{max}}] = \mu_S + \int_0^\infty \Pr\{S_{\text{max}} > \mu_S + x\} dx. \quad (37)$$

To tightly bound this integral, we select a threshold  $x_0 = \sigma_S \sqrt{2 \ln M_{\text{off}}}$  and partition the integration domain:

$$\mathbb{E}[S_{\text{max}}] \leq \mu_S + x_0 + \int_{x_0}^\infty M_{\text{off}} \exp\left(-\frac{x^2}{2\sigma_S^2}\right) dx. \quad (38)$$

By using the standard upper bound for the Q-function and evaluating the Gaussian tail integral, we obtain:

$$\mathbb{E}[S_{\text{max}}] \leq \mu_S + \sigma_S \sqrt{2 \ln M_{\text{off}}} + \frac{\sigma_S}{\sqrt{2\pi \ln M_{\text{off}}}}. \quad (39)$$

Consequently, the programming time of a single row of  $\mathbf{P}_{\text{INV}}$  can be approximated as

$$T_{\text{INV-row}} = \max\{T_{\text{INV-d}}, T_s \mathbb{E}[S_{\text{max}}]\}. \quad (40)$$

The above analysis demonstrates that the time complexity for programming time is  $O(\sqrt{\ln N_t})$ , which depends on the number of transmit antennas in MIMO systems. This result enables more accurate estimation of the programming time in massive MIMO scenarios, and the derived expressions provide closed-form guideline for assessing the programming time overhead in memristor-based implementations.

## VI. PERFORMANCE OF PROPOSED MIMO PRECODING CIRCUIT

To evaluate the performance of the proposed circuit in executing precoding algorithms, simulation experiments are conducted to analyze the impact of memristor non-idealities on BER in the downlink scenario. In the simulation, the modulation scheme is quadrature amplitude modulation (QAM) with the modulation order  $M$ .

We first evaluate the impact of  $\epsilon_p$  and  $\epsilon_{bp}$  on system BER when executing the MMSE precoding algorithm using the proposed circuit with  $N_{\text{rx}} = 16$ ,  $N_{\text{tx}} = 32$ ,  $M = 16$  and the signal-to-noise ratio (SNR) of 16 dB. Fig. 7 shows that the system BER decreases as  $\epsilon_{bp}$  increases or as  $\epsilon_p$  decreases. As shown in the dark blue region, when  $\epsilon_{bp} \geq 6$  and  $\epsilon_p \leq 3$ , the system BER remains at a low level. Under these conditions,

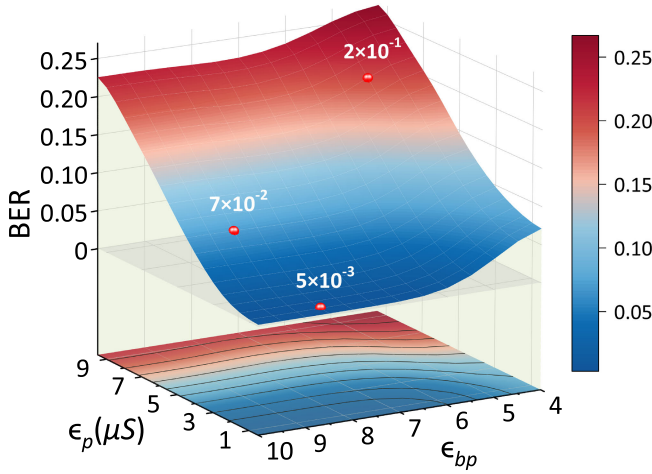


Fig. 7. Impact of  $\epsilon_p$  and  $\epsilon_{bp}$  on system BER when executing the MMSE precoding algorithm using the proposed circuit.

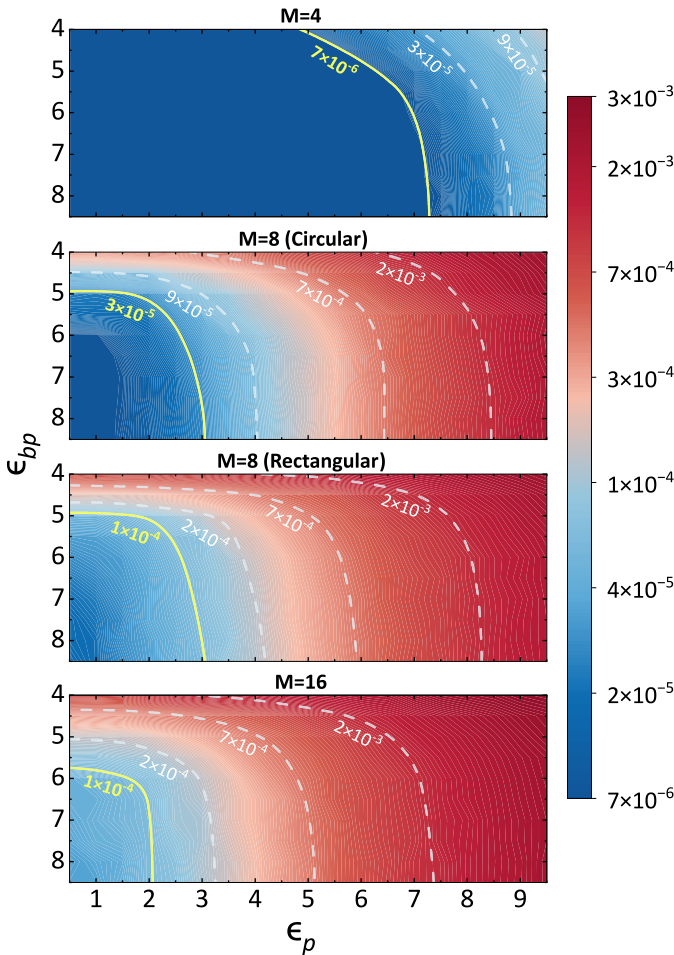


Fig. 8. BER performance of the proposed circuit under different modulation orders when executing the MMSE precoding algorithm.

the BER degradation of the proposed circuit is less than 5%, compared to FP64 precoding executed on graphics processing unit (GPU).

Next we evaluate the impact of  $\epsilon_p$  and  $\epsilon_{bp}$  with different modulation schemes on BER. The yellow contour line in Fig. 8

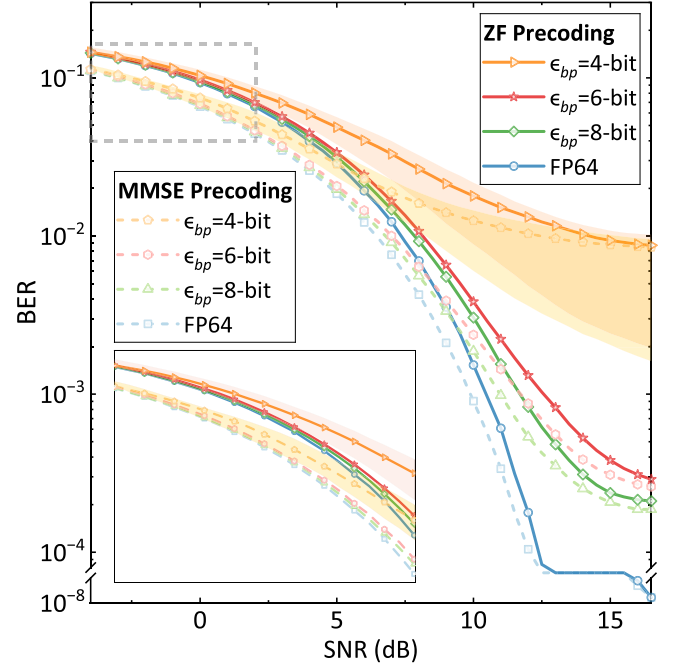


Fig. 9. BER comparison between the proposed circuit and FP64 computation under ZF and MMSE algorithms.

marks the BER obtained under FP64 computation, serving as a performance benchmark for comparison. The simulation is conducted with  $N_{\text{tx}} = 4$ ,  $N_{\text{rx}} = 32$ , and  $\text{SNR} = 10$  dB, under 4QAM, circular 8QAM, rectangular 8QAM, and 16QAM modulation schemes. As the modulation order decreases, the system BER decreases accordingly. Under 16-QAM modulation, the proposed circuit achieves a BER comparable to that of FP64 computation when  $\epsilon_{bp}$  exceeds 6 bits and  $\epsilon_p$  is below  $2 \mu\text{s}$ . Under 8-QAM modulation, comparable performance is observed with  $\epsilon_{bp}$  above 5 bits and  $\epsilon_p$  below  $3 \mu\text{s}$ . For 4-QAM modulation, the circuit achieves a comparable BER even with  $\epsilon_{bp}$  reduced to 4 bits and  $\epsilon_p$  increased to  $7 \mu\text{s}$ . These results demonstrate that reducing the modulation order significantly relaxes the precision requirements for the memristor.

In Fig. 9, we further compare the BER performance of the proposed circuit when executing ZF and MMSE algorithms under different  $\epsilon_{bp}$  values, with  $N_{\text{tx}} = 32$ ,  $N_{\text{rx}} = 16$ ,  $M = 4$  and  $\epsilon_p = 1 \mu\text{s}$ . The shaded regions in the figure represent the inter-quartile range, highlighting the spread between the 25th and 75th percentiles of the BER across Monte Carlo trials.

At low SNRs, as expected, no significant difference is observed between the proposed circuit and the FP64 computation. As the SNR increases, the BER reduction rate of the proposed circuit is significantly suppressed due to the limited bit precision of the memristor. At medium-to-high SNRs, the circuit with 4-bit memristors exhibits a significantly higher BER than that with 6-bit memristors, whereas replacing 6-bit memristors with 8-bit ones yields only marginal performance improvement.

## VII. THROUGHPUT AND ENERGY EFFICIENCY

For both ZF and MMSE precoding algorithms, the primary computational burden arises from matrix inversion, leading

TABLE I  
POWER CONSUMPTION OF MAIN COMPONENTS IN THE PROPOSED MIMO PRECODER CIRCUIT

Main Component	Key Specifications	Power Consumption	Reference
Current digital-to-analog converter (DAC)	8-bit, 0.4 ns conversion delay	1.6 mW	[18]
Voltage analog-to-digital converter (ADC)	10-bit, 0.5 ns conversion delay	41.3 $\mu$ W	[19]
Operational Amplifier (OA)	500 MHz gain-bandwidth product, 80 dB open-loop gain	12 $\mu$ W	[6]
Memristor Programming	6-bit	0.6 pJ	[20]

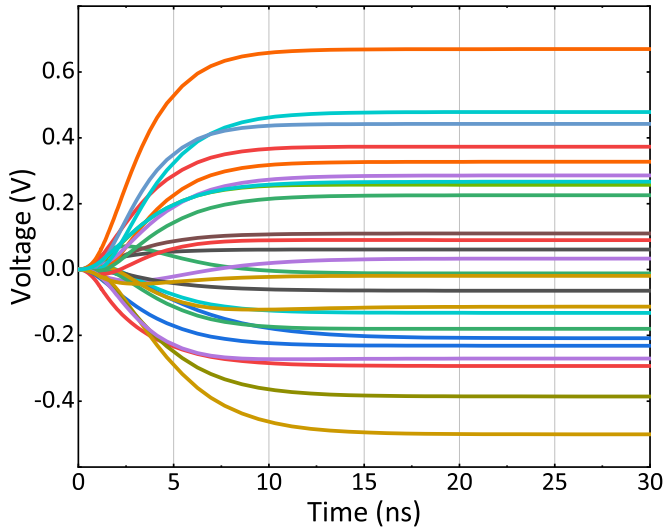


Fig. 10. Convergence time of the proposed circuit over different samples.

to a time complexity of  $O(n^3)$ , where  $n$  is the dimension of the square matrix  $\tilde{\mathbf{Z}} = \tilde{\mathbf{H}}\tilde{\mathbf{H}}^H$ . Under ideal conditions, the proposed circuit can execute ZF or MMSE precoding within a single convergence time, which remains constant regardless of the matrix size, which corresponds to a time complexity of  $O(1)$ . In the following, we provide a detailed analysis of the throughput and energy efficiency of the proposed circuit.

As an example, we determine the numbers of floating-point operations (FLOPs) required by the two precoding algorithms given  $N_{\text{rx}} = 32$  and  $N_{\text{tx}} = 16$ . Considering the operations accelerated by the proposed circuit, completing one group of symbols requires 25,600 FLOPs for ZF precoding and 25,856 FLOPs for MMSE precoding.

Studies have shown that the convergence time of the resistor crossbar is independent of the crossbar size, but depends on the factors, such as the distribution of conductance values and the eigenvalues of the correlation matrix associated with the conductance matrix [21], [22]. LTspice<sup>®</sup> transient simulation results indicate that the convergence time of the proposed circuit is less than 20 ns, as illustrated in Fig. 10, where each curve represents the temporal evolution of the output voltage vector  $\mathbf{v}_x$  in response to one group of input symbols.

Based on the estimation of programming time  $T_{\text{prog}}$  described in Section V and the single convergence time  $T_{\text{conv}}$  in Fig. 10, the total computation time  $T_{\text{total-mem}}$  can be obtained by  $T_{\text{total-mem}} = T_{\text{prog}} + T_{\text{conv}}$ . For commercial processors, the total computation time  $T_{\text{total-pro}}$  consists of data transfer time  $T_{\text{tran}}$  and floating-point computation time  $T_{\text{comp}}$ .

According to the estimation in [6], the total computation time is approximately twice the floating-point computation time  $T_{\text{total-pro}} \approx 2 \cdot T_{\text{comp}}$ .

Table I lists the energy consumption of the main components in our MIMO precoder circuit, and Table II presents the energy consumption of the four baseline general-purpose processors. The energy consumption of the proposed circuit is evaluated with reference to the method proposed in [6], focusing on its key components, while the energy consumption of the commercial processors are estimated based on their datasheets.

In addition to the energy consumption, we also quantify the area footprint of the main components in the proposed circuit. To facilitate comparison, all circuit areas are evaluated using the same 14-nm CMOS technology node adopted in [6]. Specifically, the operational amplifier occupies approximately  $50 \mu\text{m}^2$  [6], the 8-bit DAC requires  $3.07 \mu\text{m}^2$  [18], the 10-bit ADC occupies  $0.01 \text{mm}^2$  [19]. Regarding the memory elements, the area footprint of a regular RRAM cell is estimated based on the  $4F^2$  scaling assumption reported in [20]. For the proposed diagonal cell, its area is estimated by proportionally scaling the  $0.46 \mu\text{m}^2$  footprint reported for the 6T6R cell with parallel resistors in [6].

Furthermore, throughput, energy efficiency, and area efficiency are adopted as evaluation metrics to comprehensively assess the performance of the proposed circuit. They are respectively defined as

$$\text{Throughput} = \text{FLOP} / \text{Total computation time}, \quad (41)$$

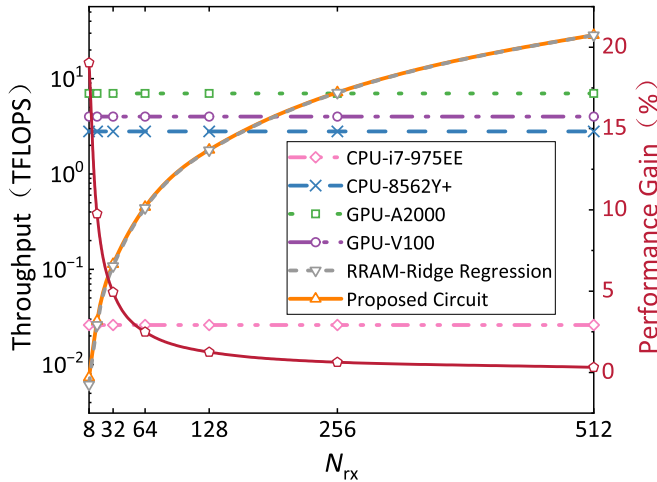
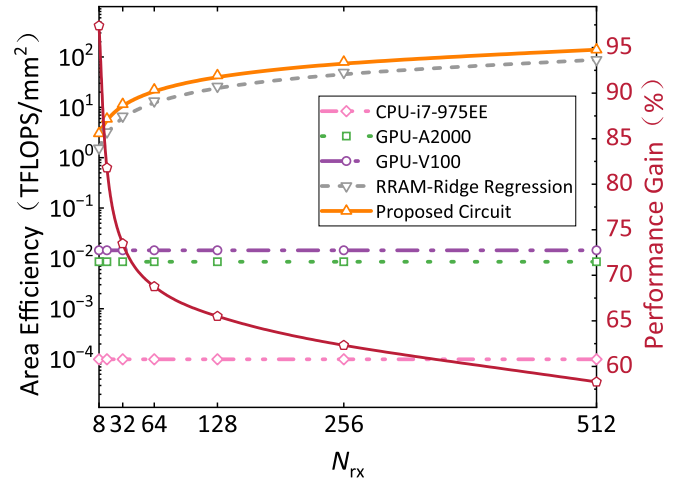
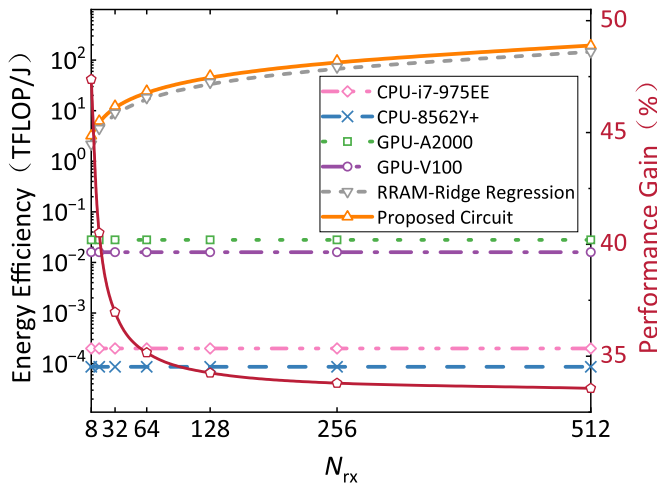
$$\text{Energy efficiency} = \text{FLOP} / \text{Energy consumption}, \quad (42)$$

$$\text{Area efficiency} = \text{Throughput} / \text{Area}. \quad (43)$$

Figs. 11-13 compare the throughput, energy efficiency, and area efficiency of the proposed circuit with those of four representative commercial processors (Intel Core i7-975EE, Intel Xeon Platinum 8562Y+, NVIDIA RTX A2000, and NVIDIA Tesla V100) as well as an analog RRAM-based baseline circuit that implements the same MMSE algorithm through ridge-regression computation [6]. The specifications of the commercial processors were obtained from their official datasheets [23], [24], [25], [26] and are summarized in Table II. These processors serve as representative general-purpose central processing units (CPUs) and GPUs, while the analog RRAM design serves as a representative in-memory-computing (IMC) baseline for the same algorithmic functionality. As can be seen from Figs. 11-13, the throughput, energy efficiency, and area efficiency of commercial processors remain nearly constant regardless of the matrix dimension, while those of the proposed memristor crossbar improve significantly as the matrix size increases.

TABLE II  
 POWER CONSUMPTION OF COMMERCIAL PROCESSORS

Processor Type	Model	Power Consumption	Peak Floating-Point Performance	Die Size	Reference
CPU	Intel Core i7-975EE	130 W	53.28 GFLOPS	263 mm <sup>2</sup>	[23]
CPU	Intel Xeon Platinum 8562Y+	300 W	5.6 TFLOPS	N/A	[24]
GPU	NVIDIA RTX A2000	70 W	8 TFLOPS	276 mm <sup>2</sup>	[25]
GPU	NVIDIA TESLA V100	250 W	14 TFLOPS	815 mm <sup>2</sup>	[26]


 Fig. 11. Throughput of the proposed circuit compared with commercial processors and an RRAM IMC baseline under varying matrix sizes ( $N_{rx} = 2N_{tx}$ ). The right-axis scale indicates the performance gain of the proposed circuit relative to the RRAM IMC baseline.

 Fig. 13. Area efficiency of the proposed circuit compared with commercial processors and an RRAM IMC baseline under varying matrix sizes ( $N_{rx} = 2N_{tx}$ ). The right-axis scale indicates the performance gain of the proposed circuit relative to the RRAM IMC baseline.

 Fig. 12. Energy efficiency of the proposed circuit compared with commercial processors and an RRAM IMC baseline under varying matrix sizes ( $N_{tx} = 2N_{rx}$ ). The right-axis scale indicates the performance gain of the proposed circuit relative to the RRAM IMC baseline.

When the matrix size  $N_{rx}$  is smaller than 16, the throughput of the proposed circuit is lower than that of the Intel Core i7-975EE. At  $N_{rx} = 256$ , the throughput becomes comparable to that of state-of-the-art commercial processors. As the matrix size increases further, the throughput of the memristor crossbar surpasses all the selected baseline processors. Compared with the RRAM IMC baseline, the proposed design

achieves a slightly higher throughput, as reflected by the performance gain curve with its scale shown on the right axis.

In terms of energy efficiency, the memristor crossbar consistently demonstrates advantages of several orders of magnitude across all matrix sizes to commercial processors. For example, the energy efficiency of the proposed circuit is 100 times higher than that of the NVIDIA RTX A2000 when  $N_{rx} = 8$ . Similar to the RRAM IMC baseline, the proposed design exhibits an increasing trend as  $N_{rx}$  grows, yet consistently achieves a 30%–50% improvement across all matrix sizes.

Regarding area efficiency, the proposed circuit outperforms the selected CPU and GPU baselines by two to three orders of magnitude. Compared with the RRAM IMC baseline, the proposed design follows the same increasing trend as the matrix size grows, while achieving 60%–100% higher efficiency across the evaluated range of  $N_{rx}$ .

In addition, although the integration of parallel resistors may raise concerns regarding potential area overhead, our design confines these resistive elements to the diagonal positions of the crossbar, resulting in a negligible footprint relative to the full crossbar. The relative area proportions among the three core components (Operational Amplifier, regular RRAM cell, proposed diagonal cell) of the proposed circuit are illustrated in Fig. 14, for small-scale crossbars (e.g.  $N_{rx} = 8$ ), the area occupied by the proposed diagonal cells is comparable to that of the regular RRAM cells, while both together constitute only a minor portion within the considered circuit components, whose overall area is largely dominated by the

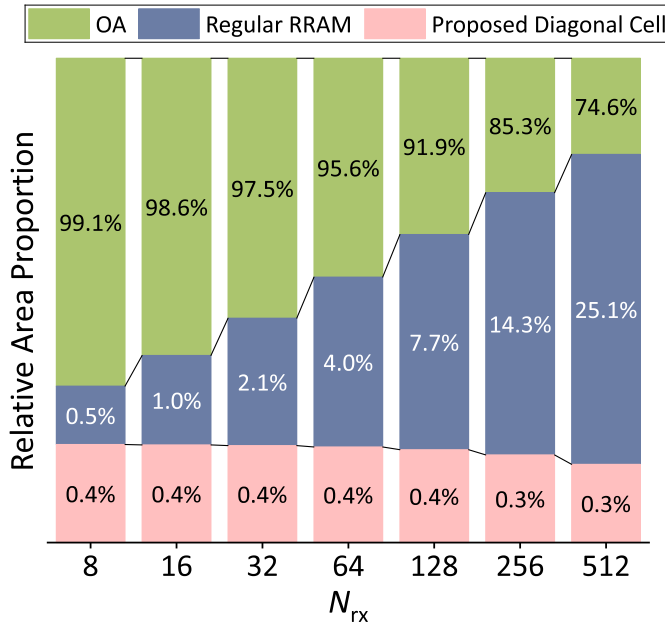


Fig. 14. Area composition of the three core components (operational amplifier, regular RRAM cell, proposed diagonal memristor cell) in the proposed circuit under varying matrix sizes ( $N_{TX} = 2N_{RX}$ ).

operational amplifiers. As the array dimension increases, the number of diagonal cells grows only linearly with  $N_{TX}$ , whereas the RRAM cells increase quadratically, leading to a clear divergence between the two. Consequently, the fraction of the diagonal cells remains nearly constant around 0.3%–0.4% and even slightly decreases, indicating that their contribution to the area of the proposed circuit becomes progressively less significant. Therefore, the integration of these additional resistors does not compromise the density advantage of the proposed IMC architecture.

## VIII. CONCLUSION

In this paper, we have proposed a memristor crossbar-based precoder circuit for accelerating linear precoding algorithm in downlink massive MIMO systems. To address the computational inaccuracies arising from the limited conductance range of memristors, we have optimized the mapping scheme that jointly considers matrix characteristics and device constraints under both Rayleigh and Kronecker channel models, reducing the relative computation error by more than 60% compared to baseline approaches. Furthermore, we have developed a probabilistic model for estimating the programming time associated with the proposed circuit. Based on this model, the closed-form upper bound and the complexity for the programming time have been obtained. Simulation results have shown that when  $\epsilon_p \leq 3$  and  $\epsilon_{bp} \geq 6$ , the BER of the proposed circuit only degrades within 5% compared to FP64 GPU. Additionally, the proposed circuit achieves high throughput while significantly enhances both energy and area efficiency in massive MIMO configurations, outperforming conventional digital processors under comparable conditions.

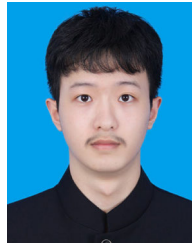
In summary, this work has identified and addressed key enabling factors for ultra-efficient precoding in massive MIMO

systems. The proposed mapping scheme and probabilistic model together bridge the gap between memristor-level non-idealities and system-level performance, laying a foundation for scalable and energy-aware baseband signal processing acceleration.

## REFERENCES

- [1] H. J. Basherlou, N. Ojaroudi Parchin, and C. H. See, "Antenna design and optimization for 5G, 6G, and IoT," *Sensors*, vol. 25, no. 5, p. 1494, Feb. 2025.
- [2] R. Kumar, D. Sinwar, and V. Singh, "Experimental evaluation of MU-MIMO in TDD environment for 5G NR using exploiting channel reciprocity," *Int. J. Inf. Technol.*, vol. 16, no. 7, pp. 4651–4666, Jun. 2024.
- [3] Z. Sun, G. Pedretti, E. Ambrosi, A. Bricalli, W. Wang, and D. Ielmini, "Solving matrix equations in one step with cross-point resistive arrays," *Proc. Nat. Acad. Sci. USA*, vol. 116, no. 10, pp. 4123–4128, Feb. 2019.
- [4] L. Xia et al., "Technological exploration of RRAM crossbar array for matrix-vector multiplication," *J. Comput. Sci. Technol.*, vol. 31, no. 1, pp. 3–19, Jan. 2016.
- [5] G. Yuan et al., "Memristor crossbar-based ultra-efficient next-generation baseband processors," in *Proc. IEEE 60th Int. Midwest Symp. Circuits Syst. (MWSCAS)*, Boston, MA, USA, Aug. 2017, pp. 1121–1124.
- [6] P. Mannonci, E. Melacarne, and D. Ielmini, "An analogue in-memory ridge regression circuit with application to massive MIMO acceleration," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 12, no. 4, pp. 952–962, Dec. 2022.
- [7] P. Zuo, Z. Sun, and R. Huang, "Extremely-fast, energy-efficient massive MIMO precoding with analog RRAM matrix computing," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 70, no. 7, pp. 2335–2339, Jul. 2023.
- [8] Y.-H. Ren, S. Yang, J.-H. Bi, and Y.-X. Zhang, "Accelerating maximum-likelihood detection in massive MIMO: A new paradigm with memristor crossbar based in-memory computing circuit," *IEEE Trans. Veh. Technol.*, vol. 73, no. 12, pp. 19745–19750, Dec. 2024.
- [9] J.-H. Bi, S. Yang, S. Chen, and P. Zhang, "High-speed ultra-energy-efficient memristor-based massive MIMO SIC detector circuit with hybrid analog-digital computing architecture," *IEEE Trans. Veh. Technol.*, vol. 74, no. 7, pp. 11495–11500, Jul. 2025.
- [10] J.-H. Bi, S. Yang, P. Zhang, and S. Chen, "Amplifier-enhanced memristive massive MIMO linear detector circuit: An ultra-energy-efficient and robust-to-conductance-error design," in *Proc. GLOBECOM IEEE Global Commun. Conf.*, Cape Town, South Africa, Dec. 2024, pp. 3968–3973.
- [11] J.-H. Bi, S. Yang, P. Zhang, and S. Chen, "In-memory massive MIMO linear detector circuit with extremely high energy efficiency and strong memristive conductance deviation robustness," in *Proc. GLOBECOM IEEE Global Commun. Conf.*, Cape Town, South Africa, Dec. 2024, pp. 728–733.
- [12] Q. Zeng et al., "Realizing in-memory baseband processing for ultra-fast and energy-efficient 6G," *IEEE Internet Things J.*, vol. 11, no. 3, pp. 5169–5183, Feb. 2024.
- [13] M. A. Albreem, A. H. A. Habbash, A. M. Abu-Hudrouss, and S. S. Ikki, "Overview of precoding techniques for massive MIMO," *IEEE Access*, vol. 9, pp. 60764–60801, 2021.
- [14] H. Lee, I. Sohn, D. Kim, and K. B. Lee, "Generalized MMSE beamforming for downlink MIMO systems," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Kyoto, Japan, Jun. 2011, pp. 1–6.
- [15] C. Li et al., "Analogue signal and image processing with large memristor crossbars," *Nature Electron.*, vol. 1, no. 1, pp. 52–59, Dec. 2017.
- [16] M. Rao et al., "Thousands of conductance levels in memristors integrated on CMOS," *Nature*, vol. 615, no. 7954, pp. 823–829, Mar. 2023.
- [17] J.-W. Jang, S. Park, G. W. Burr, H. Hwang, and Y.-H. Jeong, "Optimization of conductance change in  $\text{Pr}_{1-x}\text{Ca}_x\text{MnO}_3$ -based synaptic devices for neuromorphic systems," *IEEE Electron Device Lett.*, vol. 36, no. 5, pp. 457–459, May 2015.
- [18] S. M. I. Huq, S. Islam, N. Saqib, and S. N. Biswas, "Design of low power 8-bit DAC using PTM-LP technology," in *Proc. Int. Conf. Recent Trends Electr., Electron. Comput. Technol. (ICRTEECT)*, Telangana, India, Jul. 2017, pp. 64–69.
- [19] A. Wang and C.-J.-R. Shi, "A 10-bit 50-MS/s SAR ADC with 1 fJ/conversion in 14 nm SOI FinFET CMOS," *Integration*, vol. 62, pp. 246–257, Jun. 2018.
- [20] F. Zahoor, T. Z. A. Zulkiffi, and F. A. Khanday, "Resistive random access memory (RRAM): An overview of materials, switching mechanism, performance, multilevel cell (MLC) storage, modeling, and applications," *Nanosci. Res. Lett.*, vol. 15, no. 1, pp. 1–26, Apr. 2020.

- [21] Z. Sun, G. Pedretti, P. Mannocci, E. Ambrosi, A. Bricalli, and D. Ielmini, "Time complexity of in-memory solution of linear systems," *IEEE Trans. Electron Devices*, vol. 67, no. 7, pp. 2945–2951, Jul. 2020.
- [22] Z. Sun and R. Huang, "Time complexity of in-memory matrix-vector multiplication," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 68, no. 8, pp. 2785–2789, Aug. 2021.
- [23] Intel.(2024). *Data Sheet: I7-975 Processor Extreme Edition*. [Online]. Available: <https://ark.intel.com/content/www/us/en/ark/products/37153/intel-core-i7-975-processor-extreme-edition-8m-cache-3-33-ghz-6-40-gt-s-intel-qpi.html>
- [24] Intel.(2024). *Data Sheet: Platinum 8562Y+ Processor*. [Online]. Available: <https://www.intel.com/content/www/us/en/products/sku/237558/intel-xeon-platinum-8562y-processor-60m-cache-2-80-ghz/specifications.html>
- [25] NVIDIA.(2021). *Data Sheet: RTX A2000*. [Online]. Available: <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/rtx-a2000/nvidia-rtx-a2000-datasheet.pdf>
- [26] NVIDIA.(2017). *Data Sheet: TESLA V100*. [Online]. Available: <https://www.nvidia.com/content/dam/en-zz/zhcn/Solutions/nvidia/data-center/tesla/tesla-volta-v100-datasheet-a4-fnl-web-cn.pdf>



**Yi-Hang Ren** received the B.Eng. degree in communications engineering from Beijing University of Posts and Telecommunications (BUPT), China, in 2023. He is currently pursuing the joint Ph.D. degree in information and communication engineering with the School of Information and Communication Engineering, BUPT, and the Key Laboratory of Mathematics and Information Networks, Ministry of Education. His current research interests include baseband algorithms with in-memory computing, memristor-based circuit, and AI algorithms with in-memory computing.



**Yu-Xin Zhang** received the B.Eng. degree in communications engineering from Beijing University of Posts and Telecommunications (BUPT), China, in 2024. He is currently pursuing the joint Ph.D. degree in information and communication engineering with the School of Information and Communication Engineering, BUPT, the Key Laboratory of Mathematics and Information Networks, Ministry of Education, and the Department of IoT Technologies and Applications, China Mobile Research Institute. His research interests focus on baseband algorithms and circuit design with memristor and in-memory computing.



**Sheng Chen** (Life Fellow, IEEE) received the B.Eng. degree in control engineering from the East China Petroleum Institute, Dongying, China, in 1982, the Ph.D. degree in control engineering from the City, University of London, in 1986, and the D.Sc. degree from the University of Southampton, Southampton, U.K., in 2005. From 1986 to 1999, he held research and academic appointments with The University of Sheffield, U.K., The University of Edinburgh, U.K., and the University of Portsmouth, U.K. Since 1999, he has been with the School of

Electronics and Computer Science, University of Southampton, where he is currently a Professor of intelligent systems and signal processing. He has published over 700 research articles. He has more than 22 000 Web of Science citations with an H-index of 64 and more than 42 000 Google Scholar citations with an H-index of 87. His research interests include adaptive signal processing, wireless communications, modeling and identification of nonlinear systems, neural network and machine learning, evolutionary computation methods, and optimization. He is a fellow of the Royal Academy of Engineering, U.K., Asia-Pacific Artificial Intelligence Association (AAIA), and IET. He was one of the original ISI Highly Cited Researchers in Engineering in March 2004.



**Shaoshi Yang** (Senior Member, IEEE) received the B.Eng. degree in information engineering from Beijing University of Posts and Telecommunications (BUPT), China, in 2006, and the Ph.D. degree in electronics and electrical engineering from the University of Southampton, U.K., in 2013. From 2008 to 2009, he was a Researcher with Intel Labs, China. From 2013 to 2016, he was a Research Fellow with the School of Electronics and Computer Science, University of Southampton. From 2016 to 2018, he was a Principal Engineer with Huawei Technologies

Company Ltd., where he has made significant contributions to the products, solutions, and standardization of 5G, wideband IoT, and cloud gaming/VR. He was a Guest Researcher with the Isaac Newton Institute for Mathematical Sciences, University of Cambridge. He is currently a Full Professor with BUPT. His research interests include 5G/5G-A/6G, massive MIMO, mobile ad hoc networks, distributed artificial intelligence, and cloud gaming/VR. He is the Deputy Director of the Key Laboratory of Mathematics and Information Networks, Ministry of Education; and a Standing Committee Member of the CCF Technical Committee on Distributed Computing and Systems. He received the Dean's Award for Early Career Research Excellence from the University of Southampton in 2015, the Huawei President Award for Wireless Innovations in 2018, the IEEE TCGCC Best Journal Paper Award in 2019, the IEEE Communications Society Best Survey Paper Award in 2020, the Xiaomi Young Scholars Award in 2023, the CAI Invention and Entrepreneurship Award in 2023, the CIUR Industry-University-Research Cooperation and Innovation Award in 2023, and the First Prize of Beijing Municipal Science and Technology Advancement Award in 2023. He is an Editor of *IEEE TRANSACTIONS ON COMMUNICATIONS*, *IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY*, and *Signal Processing* (Elsevier). He was an Editor of *IEEE SYSTEMS JOURNAL* and *IEEE WIRELESS COMMUNICATIONS LETTERS*. For more details on his research progress, please refer to <https://shaoshiyang.weebly.com/>



**Ping Zhang** (Fellow, IEEE) is currently a Professor with the School of Information and Communication Engineering, Beijing University of Posts and Telecommunications; the Director of State Key Laboratory of Networking and Switching Technology; and a member of IMT-2020 (5G) Experts Panel and the Experts Panel for China's 6G Development. He served as a Chief Scientist for National Basic Research Program (973 Program), an Expert in Information Technology Division of National High-Tech Research and Development Program (863

Program), and a member of the Consultant Committee on International Cooperation of National Natural Science Foundation of China. He is an Academician of Chinese Academy of Engineering (CAE). His research interests mainly focus on wireless communications.