

Analog Array Processor with Digital Resolution Enhancement and Offset Compensation

R. Genov and G. Cauwenberghs¹
 Department of Electrical and Computer
 Engineering
 The Johns Hopkins University
 Baltimore, MD 21218, USA
 e-mail: {roman, gert}@jhu.edu

Abstract — A mixed-mode inner-product vector processor is presented. It performs high-dimensional matrix-vector multiplication on a fine-grain analog array and has a purely-digital interface. The array incorporates charge-mode analog computational cells and row-parallel analog-to-digital converters (ADC). Each of the cells includes a dynamic storage element and a charge injection device computing binary inner-product of two arguments. This eliminates constraints on memory-processor communication bandwidth leading to very high computational throughput limited only by the capacity of the array. The matrix elements are stored in the array in bit-parallel fashion, and the input vector is presented bit-serially. Digital post-processing allows to construct resulting inner-products with precision higher than that of each A/D conversion. In addition, analog computation offset errors are compensated for in a multi-chip configuration. The architecture is tailored for high-density VLSI implementation.

I. INTRODUCTION

One of the most common, but computationally most expensive operations in modern vision, image and speech processing algorithms is that of vector-matrix multiplication (VMM) in large dimensions:

$$Y^{(m)} = \sum_{n=0}^{N-1} W^{(m,n)} X^{(n)} \quad (1)$$

with N -dimensional input vector $X^{(n)}$, M -dimensional output vector $Y^{(m)}$, and $N \times M$ matrix elements $W^{(m,n)}$. In artificial neural networks, for instance, the matrix elements $W^{(m,n)}$ correspond to weights, or synapses, between neurons. The elements may also represent templates $X_n^{(m)} = W^{(m,n)}$ in a vector quantizer [1], or support vectors in a support vector machine [2].

Based on current semiconductor technologies, the only possibility to achieve real-time processing is by means of massively parallel computation. Dedicated parallel VLSI architectures have been developed to speed up VMM computation, *e.g.* [3]. The problem with most parallel systems is that they require centralized memory resources *i.e.*, RAM shared on a bus, thereby limiting the available throughput. A fine-grain, fully-parallel architecture, that integrates memory and processing elements, yields high computational throughput and high density of integration [4]. The ideal scenario (in the case of vector-matrix multiplication) is where each processor performs one multiply and locally stores one coefficient. The advantage of this is a throughput that scales linearly with the dimensions of the implemented array. The recurring problem with digital implementation is the latency in

accumulating the result over a large number of cells. Also, the extensive silicon area and power dissipation of a digital multiply-and-accumulate implementation make this approach prohibitive for very large (100-10,000) matrix dimensions.

Analog VLSI provides a natural medium to implement fully parallel computational arrays with high integration density and energy efficiency [5]. By summing charge or current on a single wire across cells in the array, low latency is intrinsic. Analog multiply-and-accumulate circuits are so small that one can be provided for each matrix element, making it feasible to implement massively parallel implementations with large matrix dimensions. Fully parallel implementation of (1) requires an $M \times N$ array of cells, each cell containing a product computing device and a storage element. Each cell (m, n) computes the product of input component $X^{(n)}$ and matrix element $W^{(m,n)}$, and dumps the resulting current or charge on a horizontal output summing line. The device storing $W^{(m,n)}$ is usually incorporated into the computational cell to avoid performance limitations due to low external memory access bandwidth. Various physical representations of inputs and matrix elements have been explored, using synchronous charge-mode [6, 7, 8, 9], asynchronous transconductance-mode [10, 11, 12], or asynchronous current-mode [13] multiply-and-accumulate circuits.

The main problem with purely analog implementation is the effect of noise and component mismatch on precision. To this end, we propose the use of hybrid analog-digital technology to simultaneously add a large number of digital values in parallel, with careful consideration of sources of imprecision in the implementation and their overall effect on the system performance. Our approach combines the computational efficiency of analog array processing with the precision of digital processing and the convenience of a programmable and reconfigurable digital interface.

A mixed-signal array architecture with binary decomposed matrix and vector elements is described in Section II. VLSI implementation is presented in Section III. Section IV quantifies the improvements obtained in system precision obtained by postprocessing the quantized outputs of the array in the digital domain. An expanded architecture using multiple processors and compensating for analog computation offset errors is discussed in Section V. Conclusions are presented in Section VI.

II. MIXED-SIGNAL ARCHITECTURE

A Internally Analog, Externally Digital Computation

The system presented is internally implemented in analog VLSI technology, but interfaces externally with the digital world. This paradigm combines the best of both worlds: it uses the efficiency of massively parallel analog computing (in particular: adding numbers in parallel on a single wire), but allows for a modular, configurable interface with other digital pre-processing and post-processing systems.

¹This work was supported by NSF MIP-9702346 and ONR N00014-99-1-0612. Chips were fabricated through the MOSIS foundry service.

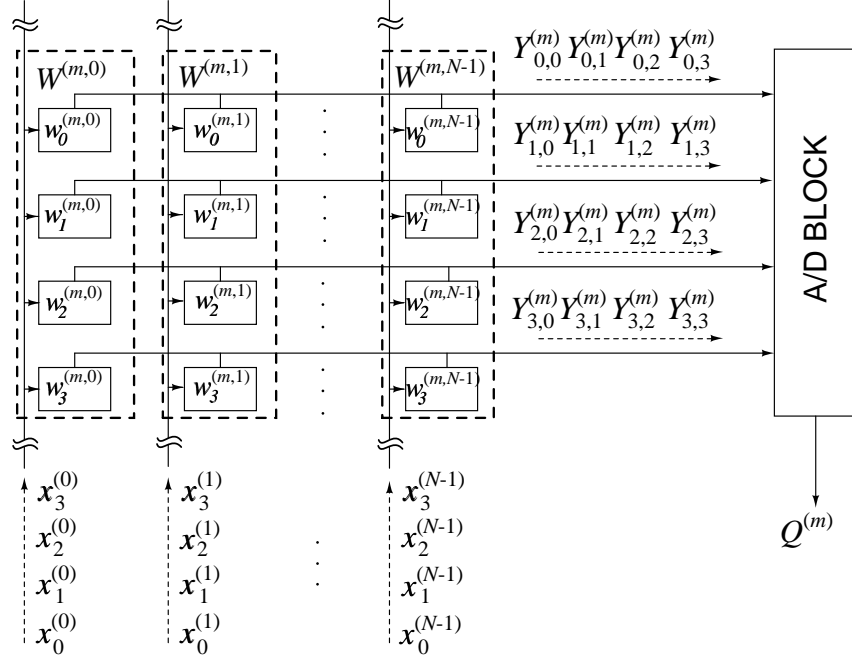


Figure 1: Block diagram of one row in the matrix with binary encoded elements $w_i^{(m,n)}$, for a single m and with $I = 4$ bits. Data flow of bit-serial inputs $x_j^{(n)}$ and corresponding partial outputs $Y_{i,j}^{(m)}$, with $J = 4$ bits.

This is necessary to make the processor a general-purpose device that can tailor the vector-matrix multiplication task to the particular application where it is being used.

The digital representation is embedded, in both bit-serial and bit-parallel fashion, in the analog array architecture (Fig. 1). Inputs are presented in bit-serial fashion, and matrix elements are stored locally in bit-parallel form. Digital-to-analog (D/A) conversion at the input interface is inherent in the bit-serial implementation, and row-parallel analog-to-digital (A/D) converters are used at the output interface.

For simplicity, an unsigned binary encoding of inputs and matrix elements is assumed here, for one-quadrant multiplication. This assumption is not essential: it has no binding effect on the architecture and can be easily extended to a standard one's complement for four-quadrant multiplication, in which the significant bits (MSB) of both arguments have a negative rather than positive weight. Assume further I -bit encoding of matrix elements, and J -bit encoding of inputs:

$$W^{(m,n)} = \sum_{i=0}^{I-1} 2^{-(i+1)} w_i^{(m,n)} \quad (2)$$

$$X^{(n)} = \sum_{j=0}^{J-1} 2^{-(j+1)} x_j^{(n)} \quad (3)$$

decomposing (1) into:

$$Y^{(m)} = \sum_{n=0}^{N-1} W^{(m,n)} X^{(n)} = \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} 2^{-(i+j+2)} Y_{i,j}^{(m)} \quad (4)$$

with binary-binary VMM partials:

$$Y_{i,j}^{(m)} = \sum_{n=0}^{N-1} w_i^{(m,n)} x_j^{(n)}. \quad (5)$$

The proposed mixed-signal approach is to compute and accumulate the binary-binary partial products (5) using an analog VMM array,

and to combine the quantized results in the digital domain according to (4).

B Array Architecture and Data Flow

To conveniently implement the partial products (5), the binary encoded matrix elements $w_i^{(m,n)}$ are stored in bit-parallel form, and the binary encoded inputs $x_j^{(n)}$ are presented in bit-serial fashion. The bit-serial format was first proposed and demonstrated in [7], with binary-analog partial products using analog matrix elements for higher density of integration. The use of binary encoded matrix elements relaxes precision requirements and simplifies storage [8].

One row of I -bit encoded matrix elements uses I rows of binary cells. Therefore, to store an $M \times N$ digital matrix $W^{(m,n)}$, an array of $MI \times N$ binary cells $w_i^{(m,n)}$ is needed. One bit of an input vector is presented each clock cycle, taking J clock cycles of partial products (5) to complete a full computational cycle (1). The input binary components $x_j^{(n)}$ are presented least significant bit (LSB) first, to facilitate the digital postprocessing to obtain (4) from (5) (as elaborated in Section IV).

Figure 1 depicts one row of matrix elements $W^{(m,n)}$ in the binary encoded architecture, comprising I rows of binary cells $w_i^{(m,n)}$, where $I = 4$ in the example shown. The data flow is illustrated for a digital input series $x_j^{(n)}$ of $J = 4$ bits, LSB first (*i.e.*, descending index j). The corresponding analog series of outputs $Y_{i,j}^{(m)}$ in (5) obtained at the horizontal summing nodes of the analog array is quantized by a bank of analog-to-digital converters (ADC), and digital postprocessing (4) of the quantized series of output vectors yields the final digital result (1).

The quantization scheme used is critical to system performance. As shown in Section IV, appropriate postprocessing in the digital domain to obtain (4) from the quantized partial products $Y_{i,j}^{(m)}$ can lead to a significant enhancement in system resolution, well beyond that of intrinsic ADC resolution. This relaxes precision requirements on the analog implementation of the partial products (5). A dense and efficient charge-mode VLSI implementation is described next.

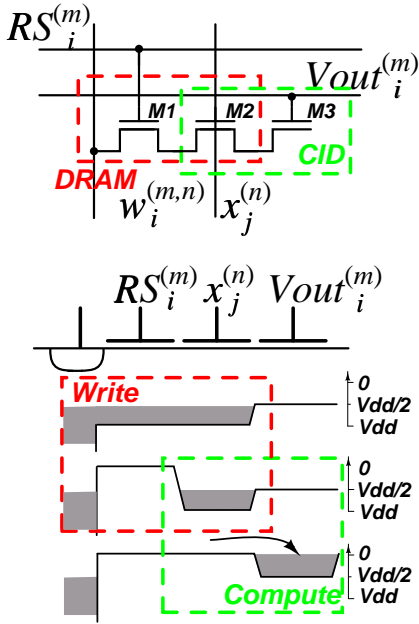


Figure 2: CID computational cell with integrated DRAM storage (top). Charge transfer diagram for active write and compute operations (bottom).

III. CHARGE-MODE VLSI IMPLEMENTATION

The elementary cell combines a CID computational unit [7, 8], computing one argument of the sum in (5), with a DRAM storage element. The cell stores one bit of a matrix element $w_i^{(m,n)}$, performs a one-quadrant binary-binary multiplication of $w_i^{(m,n)}$ and $x_j^{(n)}$, and accumulates the result across cells with common m and i indices. The circuit diagram and operation of the cell are given in Figure 2. An array of cells thus performs (unsigned) binary multiplication (5) of matrix $w_i^{(m,n)}$ and vector $x_j^{(n)}$ yielding $Y_{i,j}^{(m)}$, for values of i in parallel across the array, and values of j in sequence over time.

The cell contains three MOS transistors connected in series as depicted in Figure 2. Transistors M1 and M2 comprise a dynamic random-access memory (DRAM) cell, with switch M1 controlled by Row Select signal $RS_i^{(m)}$. When activated, the binary quantity $w_i^{(m,n)}$ is written in the form of charge stored under the gate of M2. Transistors M2 and M3 in turn comprise a charge injection device (CID), which by virtue of charge conservation moves electric charge between two potential wells in a non-destructive manner [7, 8, 14].

The cell operates in two phases: *Write* and *Compute*. When a matrix element value is being stored, $x_j^{(n)}$ is held at Vdd and $Vout$ at a voltage $Vdd/2$. To perform a write operation, either an amount of electric charge is stored under the gate of M2, if $w_i^{(m,n)}$ is low, or charge is removed, if $w_i^{(m,n)}$ is high. The amount of charge stored, ΔQ or 0, corresponds to the binary value $w_i^{(m,n)}$.

Once the charge has been stored, the switch M1 is deactivated, and the cell is ready to compute. The charge left under the gate of M2 can only be redistributed between the two CID transistors, M2 and M3. An active charge transfer from M2 to M3 can only occur if there is non-zero charge stored, and if the potential on the gate of M2 drops below that of M3 [7]. This condition implies a logical AND, *i.e.*, unsigned binary multiplication, of $w_i^{(m,n)}$ and $x_j^{(n)}$. The multiply-and-accumulate operation is then completed by capacitively sensing the amount of charge transferred onto the electrode of M3, the output summing node. To this end, the voltage on the output line, left floating after being pre-charged to $Vdd/2$, is observed. When the

charge transfer is active, the cell contributes a change in voltage

$$\Delta V_{out} = \Delta Q / C_{M3} \quad (6)$$

where C_{M3} is the total capacitance on the output line across cells. The total response is thus proportional to the number of actively transferring cells. After deactivating the input $x_j^{(n)}$, the transferred charge returns to the storage node M2. The CID computation is non-destructive and intrinsically reversible [7], and DRAM refresh is only required to counteract junction and subthreshold leakage.

The bottom diagram in Figure 2 depicts the charge transfer timing diagram for write and compute operations in the case when both $w_i^{(m,n)}$ and $x_j^{(n)}$ are of logic level 1. A logic level 0 for $w_i^{(m,n)}$ is represented as Vdd , and a logic level 1 is represented as $Vdd/2$, where Vdd is the supply voltage. For $x_j^{(n)}$, logic level 0 is represented as Vdd , and logic level 1 as GND.

Transistor-level simulation of a 512-element row indicates a dynamic range of 43 dB, and a computational cycle of 10 μs with power consumption of 50 nW per cell.

IV. QUANTIZATION AND DIGITAL RESOLUTION ENHANCEMENT

Significant improvements in precision can be obtained by exploiting the binary representation of matrix elements and vector inputs, and performing the computation (4) in the digital domain, from quantized estimates of the partial outputs (5). The effect of averaging the quantization error over a large number of quantized values of $Y_{i,j}^{(m)}$ boosts the precision of the digital estimate of $Y^{(m)}$, beyond the intrinsic resolution of the analog array and the A/D quantizers used.

A Accumulation and Quantization

The outputs $Y_{i,j}^{(m)}$ for a single m obtained from the analog array over J clock cycles can be conceived as an $I \times J$ matrix, shown in Figure 1. Elements of this matrix located along diagonals (*i.e.*, elements with a common value of $i + j$) have identical binary weight in (4). Therefore, the summation in (4) could be rearranged as:

$$Y^{(m)} = \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} 2^{-(i+j+2)} Y_{i,j}^{(m)} = \sum_{k=0}^{K-1} 2^{-(k+2)} Y_k^{(m)} \quad (7)$$

where $k = i + j$, $K = I + J - 1$ and

$$Y_k^{(m)} = \sum_{i=\kappa(k,J)}^{k-\kappa(k,I)} Y_{i,k-i}^{(m)} \quad (8)$$

with $\kappa(k, I) \equiv \max(0, k - I + 1)$ and $\kappa(k, J) \equiv \max(0, k - J + 1)$.

Several choices can be made in the representation of the signals being accumulated and quantized. One choice is whether to quantize each array output, $Y_{i,j}^{(m)}$, and accumulate the terms in (8) in the digital domain, or accumulate the terms in the analog domain and quantize the resulting $Y_k^{(m)}$. Clearly, the former leads to higher precision, while the latter has lower complexity of implementation. We opted for the former, and implemented a parallel array of low-resolution (6-bit) flash ADCs, one for each row output i .

B Row-parallel Flash A/D Conversion

Consider first the case of row-parallel flash (*i.e.*, bit-parallel) A/D conversion, where all $I \times J$ values of $Y_{i,j}^{(m)}$ are fully quantized. Figure 3 presents the corresponding architecture, shown for a single output vector component m . Each of the I horizontal summing nodes, one for each bit-plane i of component m , interfaces with a dedicated flash A/D converter producing a digital output $Q_{i,j}^{(m)}$ of

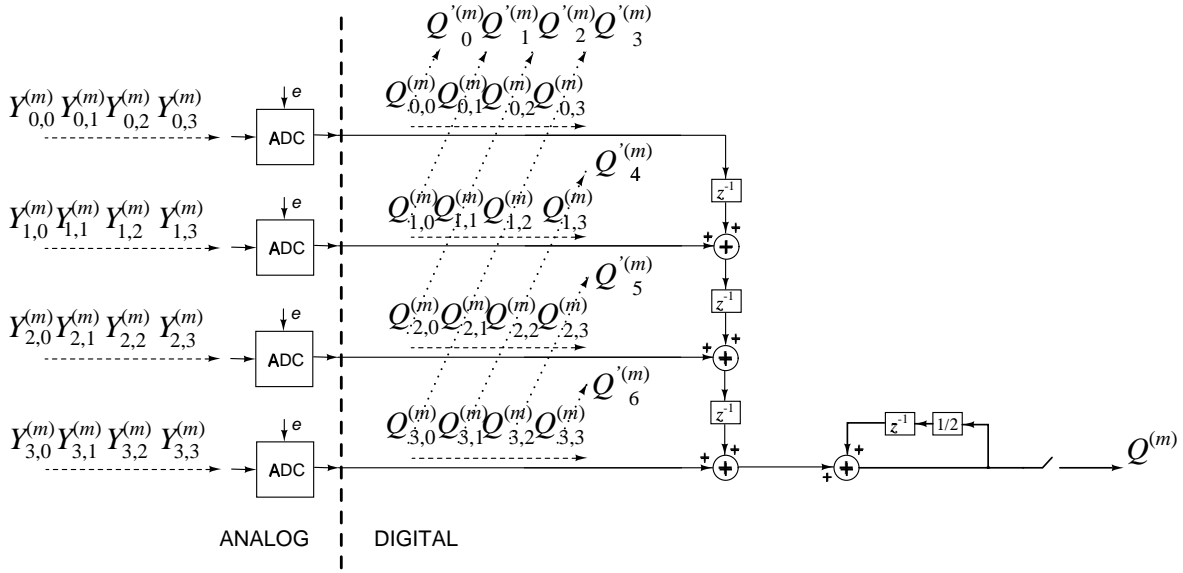


Figure 3: Diagram for the A/D quantization and digital postprocessing block in Figure 1, using row-parallel flash A/D converters. The example shown is for a single m , LSB-first bit-serial inputs, and $I = J = 4$.

L -bit resolution. The summations (8) and (7) are then performed in the digital domain:

$$Q^{(m)} = \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} 2^{-(i+j+2)} Q_{i,j}^{(m)} = \sum_{k=0}^{K-1} 2^{-(k+2)} Q_k'^{(m)}, \quad (9)$$

and

$$Q_k'^{(m)} = \sum_{i=\kappa(k,J)}^{k-\kappa(k,I)} Q_{i,k-i}^{(m)}. \quad (10)$$

A block diagram for a digital implementation is shown on the right of Figure 3, assuming LSB-first bit-serial inputs (descending index j). With radix 2, a shift-and-accumulate operation avoids the need for digital multiplication. The LSB-first bit-serial format minimizes latency and reduces the length of the register accumulating $Q^{(m)}$.

If the ADC is capable of resolving each individual binary term in the analog sum (5), then the sum is retrieved from the ADC with zero error, as if computed in the digital domain. For *zero-error* digital reconstruction, the ADC requires (at least) $N + 1$ quantization levels, that coincide with the levels of the charge transfer characteristic for any number (0 to N) of active cells along the output row of the analog array. Provided nonlinearity and noise in the analog array and the ADC are within one LSB (at the $\log_2(N + 1)$ -bit level), the *quantization error* then reduces to zero, and the output $Q^{(m)}$ is obtained at the maximum digital VMM resolution of $I + J + \log_2(N + 1)$ bits. For large arrays, this is usually more than needed, and places too stringent requirements on analog precision, $L \geq \log_2(N + 1)$.

In what follows we study the error of the digitally constructed output $Q^{(m)}$ in the practical case where the resolution of the ADC is below that of the dimensions of the array, $L < \log_2(N + 1)$. In particular, we study the properties of $Q^{(m)}$ assuming uncorrelated statistics of quantization error. The analysis yields an estimate of the gain in resolution that can be obtained relative to that of the ADC quantizers, independent of the matrix and input representation N , I , and J . The quantization is modeled as:

$$Q_{i,j}^{(m)} = Y_{i,j}^{(m)} + e_{i,j}^{(m)}, \quad (11)$$

where $e_{i,j}^{(m)}$ represents the quantization error, modeled as uniform random *i.i.d.* within one LSB. Conceptually, the error term $e_{i,j}^{(m)}$ in (11) could also include effects of noise and nonlinear distortion in the analog summation (5), although in practice the precision of the array exceeds the ADC resolution. From (9) and (11), the error in the digitally constructed output

$$Q^{(m)} = Y^{(m)} + E^{(m)}, \quad (12)$$

can then be expanded as

$$E^{(m)} = \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} 2^{-(i+j+2)} e_{i,j}^{(m)}. \quad (13)$$

Define s the full-scale range of the ADC acquiring $Q_{i,j}^{(m)}$, and S the corresponding range of the constructed digital output $Q^{(m)}$. Then according to (9),

$$S = s \sum_{i=0}^{I-1} 2^{-(i+1)} \sum_{j=0}^{J-1} 2^{-(j+1)} = s (1 - 2^{-I})(1 - 2^{-J}) \quad (14)$$

which approaches s for $I, J \rightarrow \infty$. Therefore, the full signal range is approximately equal to the output signal range of each of the ADCs.

Let the variance of the uniform quantization noise e in (11) be σ_e^2 , identical $\forall i, j$. In the *Central Limit*, the cumulative quantization error E can be roughly approximated as a normal process, with variance equal to the sum of the variances of all terms in the summation (13). Each signal component, $Q_{i,j}^{(m)}$, with quantization noise e but scaled with binary weight $2^{-(i+j+2)}$, contributes a variance $2^{-2(i+j+2)} \sigma_e^2$ in the sum (13), and the total variance σ_E^2 of the output error E is expressed as:

$$\sigma_E^2 = \sigma_e^2 \sum_{i=0}^{I-1} 2^{-2(i+1)} \sum_{j=0}^{J-1} 2^{-2(j+1)} = \sigma_e^2 \frac{1 - 2^{-2I}}{3} \frac{1 - 2^{-2J}}{3} \quad (15)$$

which approaches $(\sigma_e/3)^2$ for $I, J \rightarrow \infty$. Therefore, the signal-to-quantization-noise ratio (SQNR) approaches

$$\frac{S}{\sigma_E} \approx 3 \frac{s}{\sigma_e} \quad (16)$$

for large I and J . In other words, by quantizing each array output $Y_{i,j}^{(m)}$ instead of the combined total $Y^{(m)}$, we obtain an improvement in signal-to-quantization-noise ratio of a factor 3.

To characterize the improved precision in terms of *effective resolution* (in bits), it is necessary to relate the second order statistics of the quantization error e or E to a measure of the error indicative of resolution. There is a certain degree of arbitrariness in doing so, but in what follows we define resolution as the *median* of the absolute error *i.e.*, the (symmetric) extent of the 50 % confidence interval of the error. The choice of convention matters, because the distributions for e and E are different— e is approximately uniform, and E in the *Central Limit* is normal.

Let e be uniformly distributed in the interval $[-\Delta, \Delta]$. The median absolute value is then $\mathcal{M}_e = \frac{1}{2}\Delta$, and the variance $\sigma_e^2 = \frac{1}{3}\Delta^2$, yielding the relation

$$\mathcal{M}_e = \frac{\sqrt{3}}{2}\sigma_e \quad (17)$$

for the uniform distribution. The median absolute value for a normal distribution, in terms of the standard deviation, is approximately

$$\mathcal{M}_E = 0.675\sigma_E \quad (18)$$

This allows to express the SQNR gain in (16) as a gain in *median resolution*:

$$\frac{S}{\mathcal{M}_E} \approx 3 \frac{\sqrt{3}/2}{0.675} \frac{s}{\mathcal{M}_e} \approx 3.85 \frac{s}{\mathcal{M}_e} \quad (19)$$

or, in other words, a gain of approximately 2 bits over the resolution of each ADC.

V. MULTI-CHIP ARCHITECTURE WITH OFFSET COMPENSATION

A Multi-Chip Architecture

The proposed method of on-chip VMM computation allows for an array size of 1000×1000 , in a $0.35 \mu\text{m}$ CMOS technology implemented on a 6×6 mm die. Computations on matrices of higher dimensionality, at maximum possible precision, can be performed by using multiple VMM chips. Processors can be cascaded to expand matrix row or column spaces beyond the limits of a single chip capacity.

In the ideal case, to extend matrix row space, digital outputs of systems with shorter input vector lengths can be combined to perform a computation in a higher dimensional input space. Extension of column space is in principle also unlimited (assuming high read-out speeds). Cascading along rows of the matrix (allowing for higher dimensionality of input vectors) requires addition of digital numbers, while cascading along columns (in order to increase the number of matrix elements for a fixed input vector dimensionality) only necessitates multi-chip output multiplexing in time.

In reality, there are a number of sources of error that contribute offsets to the output of each VMM chip. These offset terms can be compensated for in the multi-chip architecture as described in the next section.

B Offset Compensation

In Section III we already considered some of the sources of computation error. They are imprecision of binary multiplications through charge sharing between potential wells in a CID unit with capacitively coupled output, and charge-mode analog addition on a single line. Because of linearity errors the result of such a computation is valid up to approximately 7 bits. We have discussed the effect of these errors and ADC resolution on the overall system precision in Sections IV.

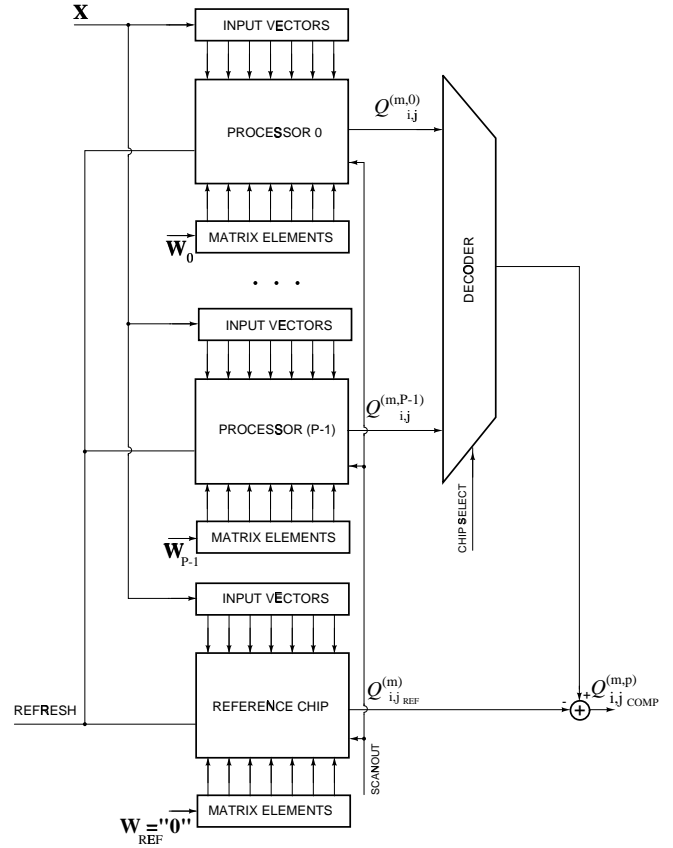


Figure 4: Multi-Chip Architecture with Offset Compensation.

Other significant sources of error in analog array-based computation are input-output feedthrough and leakage current in DRAM storage cells. The architecture presented is capable of compensating these analog computation errors in digital domain by using an extra reference VMM chip as shown in Figure 4.

The input-output feedthrough error is a result of capacitive coupling of input (vertical) lines onto the output (horizontal) lines through parasitic capacitances (metal lines overlap capacitance, gate-to-diffusion capacitance, etc). From Section III we know that the output of the analog cell ideally changes by ΔV_{out} only when both matrix element coefficient and input vector component coefficient are logic "1": $w_i^{(m,n)} = V_{dd}/2$ (charge stored) and, in the computation phase, $x_j^{(n)} = 0V$ (input switched). Any other combination of input arguments should produce zero voltage change at the output. The cell is effectively an analog AND gate. In practice, switching of the input line $x_j^{(n)}$, even when no charge is stored in CID cell, causes a small change in the output voltage, ϵ , as a result of input-output capacitive coupling. We model this effect as shown in Table 1. The output of a single cell is denoted as $y_{i,j}^{(m,n)}$.

From equations (5) and (11), taking into account the input-output feedthrough error, the VMM quantized output partials of the p -th processor in a $(P - 1)$ -processor architecture can now be expressed as:

$$\begin{aligned} Q_{i,j}^{(m,p)} &= \sum_{n=0}^{N-1} y_{i,j}^{(m,n,p)} + e_{i,j}^{(m,p)} \\ &= (1 + \epsilon) \sum_{n=0}^{N-1} w_i^{(m,n,p)} x_j^{(n)} \end{aligned}$$

Table 1: Input-output mapping of CID/DRAM computational cell. The input-output feedthrough error is introduced when the input is logic “1”.

$x_j^{(n)}$	$w_i^{(m,n)}$	$y_{i,j}^{(m,n)}$
0	0	0
0	1	0
1	0	ϵ
1	1	$1+\epsilon$

$$\begin{aligned}
& + \epsilon \sum_{n=0}^{N-1} (1 - w_i^{(m,n,p)}) x_j^{(n)} + e_{i,j}^{(m,p)} \\
& = \sum_{n=0}^{N-1} w_i^{(m,n,p)} x_j^{(n)} + \epsilon \sum_{n=0}^{N-1} x_j^{(n)} + e_{i,j}^{(m,p)}. \quad (20)
\end{aligned}$$

The term $\epsilon \sum_{n=0}^{N-1} x_j^{(n)}$ in equation (20) is an offset term proportional to the number of active vector components. To compensate for this term an identical extra, reference, chip is used in the multi-chip architecture as shown in Figure 4. The same inputs are presented to all chips in the system while only logic “0” matrix element coefficients are stored in the reference chip. It outputs quantized partials of the form:

$$Q_{i,j}^{(m)}{}_{REF} = \epsilon \sum_{n=0}^{N-1} x_j^{(n)} + e_{i,j}^{(m)}{}_{REF}. \quad (21)$$

In order to compensate for feedthrough offsets, once computation has been performed on chips, the output of the reference chip is subtracted from the output of processors in digital domain:

$$\begin{aligned}
Q_{i,j}^{(m,p)}{}_{COMP} &= Q_{i,j}^{(m,p)} - Q_{i,j}^{(m)}{}_{REF} \\
&= \sum_{n=0}^{N-1} w_i^{(m,n,p)} x_j^{(n)} + e_{i,j}^{\prime(m,p)}, \quad (22)
\end{aligned}$$

where

$$e_{i,j}^{\prime(m,p)} = e_{i,j}^{(m,p)} - e_{i,j}^{(m)}{}_{REF}. \quad (23)$$

Another important source of errors requiring specific consideration is DRAM leakage current. In a standard two-level DRAM cell, the exact amount of charge stored is not crucial. It is used only for binary operations of readout and refresh, and a slow refresh scheme can be used. In contrast, the CID/DRAM cell produces an analog output which is proportional to the amount of charge stored in the charge injection device (6). Over multiple computation cycles, different rows are being refreshed, producing differences in the temporal decay profile of charge stored along rows of cells.

In the multi-chip configuration shown in Figure 4, the refresh clock of the same frequency as in standard DRAMs is used. It is fed to all processors, including the reference chip. Using a reference chip with all cells containing logic “0” coefficients and refreshed synchronously with processor chips ensures the same charge decay profile in its cells and voltage increase profile at its outputs. To compensate for charge decays caused by leakage current, the outputs from the reference chip are subtracted from the outputs of the corresponding rows of each of the processors.

Therefore, both input-output feedthrough input-dependent offset and charge decay input and time-dependent offset are compensated for in the multi-chip VMM architecture by using one reference chip

supplied with identical inputs, synchronous refresh clock and all logic “0” matrix elements. Subtraction of outputs of equivalent rows in digital domain eliminates both input-dependent and temporal errors.

VI. CONCLUSIONS

A charge-mode VLSI architecture for parallel vector-matrix multiplication in large dimensions ($N, M = 100-10,000$) has been presented. An internally analog, externally digital architecture offers the best of both worlds: the density and energetic efficiency of an analog VLSI array, and the noise-robustness and versatility of a digital interface. The combination of analog array processing and digital post-processing also enhances the precision of the digital VMM output, exceeding the resolution of the quantized analog array outputs by 2 bits. A reference subtraction scheme with one additional processor also compensates for clock feedthrough and charge leakage in the array.

Fine-grain massive parallelism and distributed memory, in an array of 3-transistor CID/DRAM cells, provides a computational efficiency (bandwidth to power consumption ratio) exceeding that of digital multiprocessors and DSPs by several orders of magnitude.

REFERENCES

- [1] A. Gersho and R.M. Gray, *Vector Quantization and Signal Compression*, Norwell, MA: Kluwer, 1992.
- [2] V. Vapnik, *The Nature of Statistical Learning Theory*, 2nd ed., Springer-Verlag, 1999.
- [3] J. Wawrzyniec, et al., “SPERT-II: A Vector Microprocessor System and its Application to Large Problems in Backpropagation Training,” in *Advances in Neural Information Processing Systems*, Cambridge, MA: MIT Press, vol. 8, pp 619-625, 1996.
- [4] J.C. Gealow and C.G. Sodini, “A Pixel-Parallel Image Processor Using Logic Pitch-Matched to Dynamic Memory,” *IEEE J. Solid-State Circuits*, vol. **34**, pp 831-839, 1999.
- [5] A. Kramer, “Array-based analog computation,” *IEEE Micro*, vol. **16** (5), pp. 40-49, 1996.
- [6] A. Chiang, “A programmable CCD signal processor,” *IEEE Journal of Solid-State Circuits*, vol. **25** (6), pp. 1510-1517, 1990.
- [7] C. Neugebauer and A. Yariv, “A Parallel Analog CCD/CMOS Neural Network IC,” Proc. IEEE Int. Joint Conference on Neural Networks (IJCNN’91), Seattle, WA, vol. **1**, pp 447-451, 1991.
- [8] V. Pedroni, A. Agrat, C. Neugebauer, A. Yariv, “Pattern matching and parallel processing with CCD technology,” Proc. IEEE Int. Joint Conference on Neural Networks (IJCNN’92), vol. **3**, pp 620-623, 1992.
- [9] G. Han, E. Sanchez-Sinencio, “A general purpose neuro-image processor architecture,” Proc. of IEEE Int. Symp. on Circuits and Systems (ISCAS’96), vol. **3**, pp 495-498, 1996.
- [10] M. Holler, S. Tam, H. Castro and R. Benson, “An Electrically Trainable Artificial Neural Network (ETANN) with 10,240 Floating Gate Synapses,” in *Proc. Int. Joint Conf. Neural Networks*, Washington DC, pp 191-196, 1989.
- [11] F. Kub, K. Moon, I. Mack, F. Long, “Programmable analog vector-matrix multipliers,” *IEEE Journal of Solid-State Circuits*, vol. **25** (1), pp. 207-214, 1990.
- [12] G. Cauwenberghs, C.F. Neugebauer and A. Yariv, “Analysis and Verification of an Analog VLSI Incremental Outer-Product Learning System,” *IEEE Trans. Neural Networks*, vol. **3** (3), pp. 488-497, May 1992.
- [13] A.G. Andreou, K.A. Boahen, P.O. Pouliquen, A. Pavasovic, R.E. Jenkins, and K. Strohhahn, “Current-Mode Subthreshold MOS Circuits for Analog VLSI Neural Systems,” *IEEE Transactions on Neural Networks*, vol. **2** (2), pp 205-213, 1991.
- [14] M. Howes, D. Morgan, Eds., *Charge-Coupled Devices and Systems*, John Wiley & Sons, 1979.