

TD-dAJC: A 2pJ/pixel Time-Domain Weight and Integrating-MAC based direct-Analog-to-MJPEG Compression for Video Sensor Nodes

Gourab Barik*, Harshit Naman*, Yudhajit Ray, Shreyas Sen

Purdue University, USA

*Equally Credited Authors

The increase in IoT applications has made cameras ubiquitous in various fields spanning from healthcare and autonomous vehicles to energy-constrained applications like battery-powered wearables. This proliferation has created a need for low-power camera modules for resource-constrained IoT nodes/edge devices. As the quality of images improves, these cameras produce a data deluge. Standard digital cameras generate data in the order of ~Gbps, which requires compressing the data within the camera module (e.g., SONY IMX 317) using lossy compression schemes like MJPEG or H.264. For instance, a 4K 12fps video generates around ~1 Gbps of data, which can be compressed to ~50 Mbps using MJPEG. Compared to H.264, MJPEG has lower latency and lower complexity, making it suitable for real-time edge applications. Standard digital data compression suffers from 1) high ADC energy consumption digitizing all samples, 2) Intermediate data storage requirements, and 3) huge computation power to compress the data, making it less desirable for edge cameras. The emergence of low-power camera imagers [1], such as the STM VD55G1 with power consumption as low as 3-10mW, has increased the demand for efficient compression techniques that minimize computational power as shown in Fig.1(a) and (b).

Recent advancements in Time domain [2]-[7] and analog computation [8]-[14], particularly in AI-hardware and In-sensor analytics, have motivated a shift towards performing computations before digitization. This approach not only leverages the low-power analog computing but also reduces ADC conversion energy, as only a fraction of samples needs to be digitized, while eliminating the need for large intermediate storage units. In Fig.1(c), [8], a switched-capacitor (SC) MAC(Multiply-accumulate) approach has been demonstrated for direct analog-to-MJPEG (dAJC) compression, offering better energy efficiency than its digital counterparts. However, it suffers from process-dependent SC-Discrete Cosine Transform (DCT) weights, power-hungry intermediate buffers, larger area, and speed limitations. In this work as shown in Fig.1(d) and (e), those concerns are addressed through a mixed-signal computing macro that utilizes Time-Domain weight and Integration for MAC and division to enable dAJC (TD-dAJC). The TD-dAJC technique maps DCT weights to time pulses extracted from a 32-Phase PLL(by XORing the selected phases), making the weights PVT invariant. PVT-tolerance of the integrating MAC is achieved using built-in current control, enhancing computed MAC accuracy. Unlike SCs, the TD technique eliminates need for large capacitances and hence intermediate power-hungry buffers, making the solution significantly lower area and lower power, while achieving 20x faster operation speeds (5MSps [8] to 100MSps). Time pulse-based computing makes TD-dAJC more suitable for lower-scaled technology nodes compared to voltage or current-based analog techniques. This work achieves a 25x improvement in energy per pixel compared to digital implementations and a 13x improvement over previous SC-based implementation with an increased frame size and frame rate as shown in Fig.1(f).

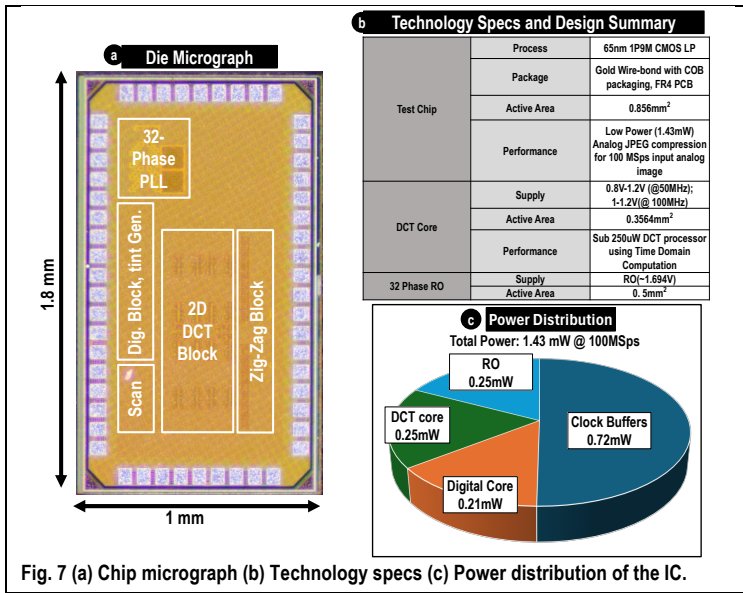
The unit TD-MAC circuit is shown in Fig. 2(a). A source-degenerated native NMOS (M0) is used for linear voltage-to current (V-I) conversion, with $V_{th} \sim 0$ mV, allowing an input swing of $(V_{dd}/2)$ which helps maintain a high SNR as the signal traverses through MAC stages. The output swing is matched to input swing (unity gain), eliminating the need for buffers. The resistor(R) is adjusted to compensate for process variations in the V-I conversion ratio and capacitor value. The input voltage signal is converted from V-I using M0 and mirrored using cascoded current mirror and is integrated on a capacitance CMAC, for a specific time proportional to the DCT weight. Switches S1 and S0 are closed during the DCT weight time pulse, leading to multiplication. M3 acts as a cascode device making the integration current less sensitive to (VOUT). After multiplication, S0 and S1 open, reducing static current during idle phases. This process repeats for 8 clock cycles, accumulating 8 multiplications for 8 input samples and corresponding TD weights onto CMAC, realizing an 8x1 MAC. Fig.2(a) (bottom) shows the Pseudo-Differential MAC units (PD-MAC) for handling signed weights. Switches IPxDP/N adjust connections based on the DCT weight's sign. For positive weights, the input (Vin+) goes to VINP, while VINN connects to a DC common mode (CM) voltage (Vin-); for negative weights, this is reversed. Fig.2(b) shows the time domain MAC computation waveforms. For negative weights t3 and t4, V3 and V4 connect to VINN, with VINP connected to the CM voltage. The VOUT waveform shows subtraction for these weights. Fig.2(c) shows the system diagram with MJPEG blocks. The MJPEG compression applies the 2D-DCT to 8x8 pixel blocks. After DCT, the computed 8x8 values are element-wise divided by a quantization matrix (Q50), which controls compression ratio and image quality. Following quantization, the compressed

values are reordered in a zig-zag pattern to group lower-frequency components. Fig.2(d) shows TD weights generation from the 32 phases of a 100MHz clock. Selected phases are used by an XOR gate-based weight generator to create 7-unique time pulses corresponding to each 7-unique value in DCT weight matrix. These pulses are then directed to the MAC units for the DCT and Quantization operation. values are element-wise divided by a quantization matrix (Q50), which controls compression ratio and image quality. Following quantization, the compressed values are reordered in a zig-zag pattern to group lower-frequency components. Fig.2(d) shows TD weights generation from the 32 phases of a 100MHz clock. Selected phases are used by an XOR gate-based weight generator to create 7-unique time pulses corresponding to each 7-unique value in DCT weight matrix. These pulses are then directed to the MAC units for the DCT and Quantization operation.

Fig.3 shows the circuit-diagram of the TD-dAJC, the input image processing setup and the post-processing. The input image is divided into 8x8 pixel matrices, with pixel values ranging from 0-255, mapped to the input voltage range. These voltages are sent to the Vin+ port of the IC using an Arbitrary Waveform Generator (AWG) as shown in Fig.3(a). The input sampling clock (CLK_AWG) is provided to the IC's clocking circuitry. The synchronizer block selects appropriate clock phases from the 32-phase clock (for Clock and Data Alignment) and sends them to the TD-DCT weight generation block. This block generates the input sampling clocks, DCT weight time pulses and Quantization(Q) time pulses, which are then buffered and sent to the MAC units in the DCT computation core, and the Q-block as shown in Fig.3(b). The DCT circuit operates in 2-stages, processing an 8x8 matrix of pixels at a time as shown in Fig 3(c): the first stage(1D-DCT) consists of 8-PD-MAC units. These units receive sampling clocks (IP1DP/N) and DCT weight pulses (tint1D), computing MAC output for one column in 8-cycles at the output capacitor. In the 9th cycle, these outputs are sampled to the next stage using 2D-DCT sampling clocks (IP2DP/N). The second stage(2D-DCT) comprises eight such 1D-DCT units, with DCT weight pulses (tint2D) provided at an 8x lower frequency. The DCT core computes the 2D-DCT of an 8x8 image block in 72 cycles (9 cycles x 8 columns), resetting the 1D-DCT block after each column. Following 2D-DCT, Quantization occurs using in-situ Q-blocks during the off-phase of the 72nd cycle. This block performs elementwise TD-division by sampling the 2D-DCT output to the input of a TD-MAC unit cell that subtracts an appropriate ratio of the output from the same 2D-DCT output capacitor based on the Q50 matrix as shown in Fig.3(e). This technique reduces area by eliminating the need for an extra capacitor by reusing the same 2D-DCT integration capacitors. The quantized samples are serialized such that significant samples are first followed by insignificant samples using a zig-zag traversal block with an analog mux and a digital controller, which are then sent to the output. The output samples, when sent to an off-chip comparator, activates the ADC only for significant samples (typically 5-10%), while the Run-Length-Encoding (RLE) is activated for insignificant samples, saving ADC energy by digitizing only significant samples as shown in Fig.3(d).

The source degeneration resistors in the MAC units use a 5-bit binary controllability to compensate for process variations. Fig.4(a) illustrates that the integration current can be adjusted to restore the integrated values to nominal levels in both fast and slow process corners. The rms noise current contributions of MOS and resistor in the unit TD-MAC are used to determine the integrated rms noise voltage at the 2D-DCT output after 72 cycles, which is ~20 μ V for an output capacitor of 200 fF. For an 8-bit ADC, the rms noise voltage must be below 50 μ V to achieve a target SNR of 16 dB. Thus, the design meets the noise requirements. Utilizing a larger output capacitor (~200 fF) aid in achieving lower rms noise voltage as shown in Fig.4(b). The capacitor sizing helps handle non-idealities(clock feedthrough, charge injection) and (switching and KT/C) noise. In the measurement setup (Fig.4(c)), a 4K RGB image is decomposed into 3-grayscale channels. The pixels are serialized to analog voltages and sent into the IC using an AWG. The compressed samples are captured and exported via an oscilloscope for analysis. The reconstructed image is generated by applying inverse-DCT on the measured analog samples in MATLAB and compared against the original image to verify compression fidelity. The results in Fig.4(d) show PSNR values of 30.1 dB for the MNIST dataset and 29.7 dB for 4K-checkerboard image.

Fig.5(a) highlights measurement results showing power consumption vs. VDD plot. The PSNR variation vs threshold voltage in Fig.5(b), indicates the measured optimum threshold value (V_{th}) for the ADC for best PSNR. Fig.5(c) shows input data of a 4K image, computation duration, and 100 MSps compressed analog samples. A zoomed in batch of 64 analog samples is also presented. The shmoo plot (Fig.5(d)) shows the operability of the ASIC for different supply voltages at 100 MSps. Fig.6 compares the test chip performance with state-of-the-art JPEG encoders [15]-[17] and DCT compression circuits [18]-[19] exhibiting a data rate of 100 MSps (20x improvement), enabling a resolution of 4K 12fps/HD 30fps RGB, and an energy efficiency of only 2pJ/pixel (25x improvement). Fig. 7 shows the die micrograph, technology specs, and the power distribution of the IC.



Acknowledgements

This work was supported by Quasistatics Inc.

References

- [1] I. Park, W. Jo *et al*, "A 640 × 640 Fully Dynamic CMOS Image Sensor for Always-On Operation", JSSC2020. <https://doi.org/10.1109/JSSC.2019.2959486>.
- [2] P. -C. Wu *et al.*, "A 28nm 1Mb Time-Domain Computing-in-Memory 6T-SRAM Macro with a 6.6ns Latency, 1241GOPS and 37.01TOPS/W for 8b-MAC Operations for Edge-AI Devices," ISSCC 2022. <https://doi.org/10.1109/ISSCC42614.2022.9731681>.
- [3] A. Sayal, *et al*, "14.4 All-Digital Time-Domain CNN Engine Using Bidirectional Memory Delay Lines for Energy-Efficient Edge Computing," ISSCC 2019, San Francisco, CA, USA, 2019, pp. 228-230. <https://doi.org/10.1109/ISSCC.2019.8662510>.
- [4] F. Chen, *et al*, "A 108nW 0.8mm² Analog Voice Activity Detector (VAD) Featuring a Time-Domain CNN as a Programmable Feature Extractor and a Sparsity-Aware Computational Scheme in 28nm CMOS," ISSCC 2022. <https://doi.org/10.1109/ISSCC42614.2022.9731720>.
- [5] Z. Chen *et al*, "A Time-Domain Computing Accelerated Image Recognition Processor with Efficient Time Encoding and Non-Linear Logic Operation," JSSC 2019. <https://doi.org/10.1109/JSSC.2018.2883394>.
- [6] D. Miyashita, *et al*, "A Neuromorphic Chip Optimized for Deep Learning and CMOS Technology With Time-Domain Analog and Digital Mixed-Signal Processing," JSSC 2019. <https://doi.org/10.1109/JSSC.2017.2712626>.
- [7] J. Lou, *et al*, "All-Digital Time-Domain Compute-in-Memory Engine for Binary Neural Networks With 1.05 POPS/W Energy Efficiency," ESSCIRC 2022. <https://doi.org/10.1109/ESSCIRC55480.2022.9911382>.
- [8] K. Gaurav Kumar, *et al*, "A 65 nm 2.02 mW 50 Mbps Direct Analog to MJPEG Converter for Video Sensor Nodes using low-noise Switched Capacitor MAC-Quantizer with automatic calibration and Sparsity-aware ADC," CICC 2023. <https://doi.org/10.1109/CICC57935.2023.10121318>.
- [9] Miscuglio, M., Gui, Y., Ma, X. *et al*. Approximate analog computing with memtronic circuits. *Commun Phys* 4, 196 (2021). <https://doi.org/10.1038/s42005-021-00683-4>.
- [10] J. Liang, *et al*, "An Offset-Cancelling Discrete-Time Analog Computer for Solving 1-D Wave Equations," JSSC 2021. <https://doi.org/10.1109/JSSC.2021.3074003>.
- [11] J. Song, *et al*, "A 16Kb Transpose 6T SRAM In-Memory-Computing Macro based on Robust Charge-Domain Computing," ASSC 2021. <https://doi.org/10.1109/A-SSCC53895.2021.9634747>.
- [12] E. H. Lee and S. S. Wong, "Analysis and Design of a Passive Switched-Capacitor Matrix Multiplier for Approximate Computing," JSSC 2017. <https://doi.org/10.1109/JSSC.2016.2599536>.
- [13] B. Zhang *et al.*, "A 177 TOPS/W, Capacitor-based In-Memory Computing SRAM Macro with Stepwise-Charging/Discharging DACs and Sparsity-Optimized Bitcells for 4-Bit Deep Convolutional Neural Networks," CICC 2022. <https://doi.org/10.1109/CICC53496.2022.9772781>.
- [14] Z. Chen, *et al*, "DCT-RAM: A Driver-Free Process-In-Memory 8T SRAM Macro with Multi-Bit Charge-Domain Computation and Time-Domain Quantization," CICC 2022. <https://doi.org/10.1109/CICC53496.2022.9772782>.

- [15] N. Reynders *et al*, "27.3 A 210mV 5MHz variation-resilient near-threshold JPEG encoder in 40nm CMOS," ISSCC 2014. <https://doi.org/10.1109/ISSCC.2014.6757511>.
- [16] Yu Pu, *et al*, "An ultra-low-energy/frame multi-standard JPEG co-processor in 65nm CMOS with sub/near-threshold power supply," ISSCC 2009. <https://doi.org/10.1109/ISSCC.2009.4977350>.
- [17] S. Kawahito *et al.*, "A compressed digital output CMOS image sensor with analog 2-D DCT processors and ADC/quantizer," ISSCC 1997. <https://doi.org/10.1109/ISSCC.1997.585326>.
- [18] M. Pankaala, *et al*, "An Analog 2-D DCT Processor," in TCAS Video 2006. <https://doi.org/10.1109/TCSVT.2006.882392>.
- [19] I. Park, *et al*, "A 640 × 640 Fully Dynamic CMOS Image Sensor for Always-On Operation," JSSC 2020. <https://doi.org/10.1109/JSSC.2019.2959486>.