

D3TA: 38.9TOPS/W Transformer Accelerator with Dual-Port 3T-eDRAM Digital Compute-In-Memory using HyperAttention and Triple-Sparsity-Handling

Donghyuk Kim*, In-Jun Jung*, Geonwoo Ko, Seri Ham, Cuong Duong Manh, Gyeongrok Yang, and Joo-Young Kim
School of Electrical Engineering,

Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Republic of Korea
{kar02040, injun, geonwooko, seri1215, cuongdm1410, toddlerf, jooyoung1203}@kaist.ac.kr

*Equal Contribution

Abstract—This paper presents D3TA that accelerates transformers with high system performance. It introduces three key features: 1) HyperAttention (HA) that supports maximum pipelining in memory to minimize the external memory access (EMA) and increase CIM utilization, 2) triple sparsity handling CIM engine that skips redundant near-zero values for energy and latency reduction, and 3) high throughput Dual-port 3T-eDRAM Digital CIM (D3CIM) macro with charge-recycling (CR) and bit-parallel dual radix-4 booth (BPDB) MAC tree group to increase energy- and area-efficiency. The D3TA achieves $2.6\times$ latency reduction in BERT-Base with $2.0\times$ EMA reduction and 1.28 -to- $9.75\times$ higher system FoM than previous works.

Index Terms—Transformer, Dataflow, Sparsity, 3T-eDRAM cell, Computing-in-Memory

I. INTRODUCTION

Transformers demonstrate remarkable performance in AI domains, and digital CIM based transformer accelerators [1]–[3] have emerged to address memory bandwidth (BW) bottlenecks. However, their excessive multi-head attention (MHA) and matrix multiplication kernels cause three major system-level challenges, as shown in Fig. 1. First, as token length increases, MHA exponentially generates intermediate data of Q, K, V, score (S), and probability matrices (P), increasing EMA and lowering CIM utilization. Second, transformers perform redundant compute due to numerous near-zeros in activations (A) and weights (W) [4], as well as weakly related tokens for output (O), increasing inefficiency in CIM macros. Third, achieving high throughput in CIM macros incurs substantial energy consumption and area overhead due to multi-row activation, multi-bit-line switching, and large bit-parallel MAC. We propose D3TA, D3CIM based Transformer Accelerator equipped with three key features: 1) HA that reduces EMA and enhances CIM utilization, 2) CIM engine (CE) with triple sparsity handler (TSH) for A/W/O sparsity, and 3) D3CIM macro supporting CR and BPDB MAC tree group to reduce energy and area overhead.

II. D3TA ARCHITECTURE

Fig. 2 shows the overall D3TA architecture, consisting of 2 HA clusters (HACs). Each HAC includes 8 CEs, a streaming vector processing unit (SVPU), an 88KB attention

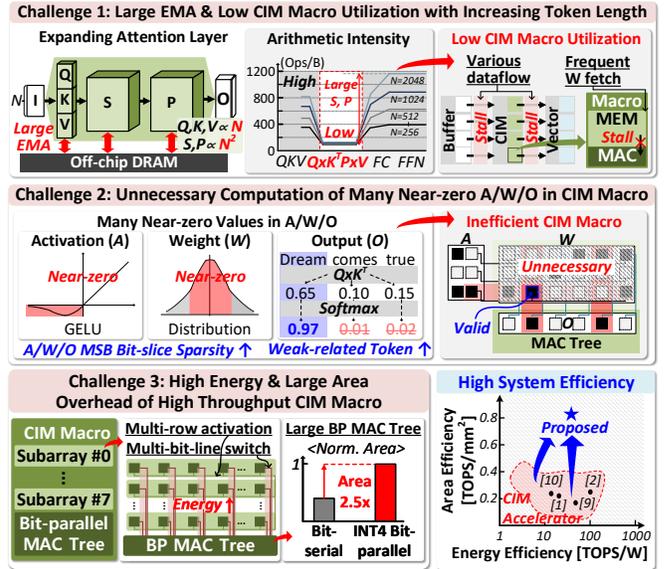


Fig. 1. System-level challenges for transformer accelerator

buffer, a reconfigurable HAC network (RHCN) for efficient data movement, and an HAC controller. A CE consists of a CE controller, input buffer, D3CIM with 544×320 cells, D3CIM accumulator for partial sum (PSUM), PSUM buffer for output of D3CIM, and TSH. A TSH consists of a triple sparsity controller (TSC), early zero-gating (EZG) bit-mask (BM) generator, and triple BM buffer for managing triple sparsity. The SVPU supports inter-CE aggregation (ICEA) for 8 CEs, quantization for 2’s complement (2C) and sign-magnitude (SM), softmax, GELU, transpose, and compressor for activations.

A. HyperAttention (HA)

Fig. 3 illustrates HA dataflow that fully retains Q, K, V, S, and P on chip using a tile-basis operation via attention buffer, and pipelines QKV generation and MHA with two optimized computations. Q, K, and V are generated on a per-head basis, then a group of S, P, and attention output rows are computed. Furthermore, RHCN and in-memory pipeline (IMP) provide

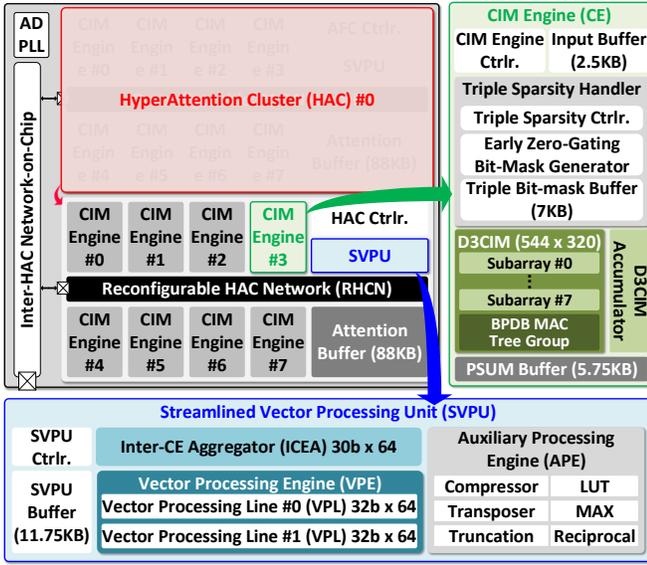


Fig. 2. Overall architecture of D3TA

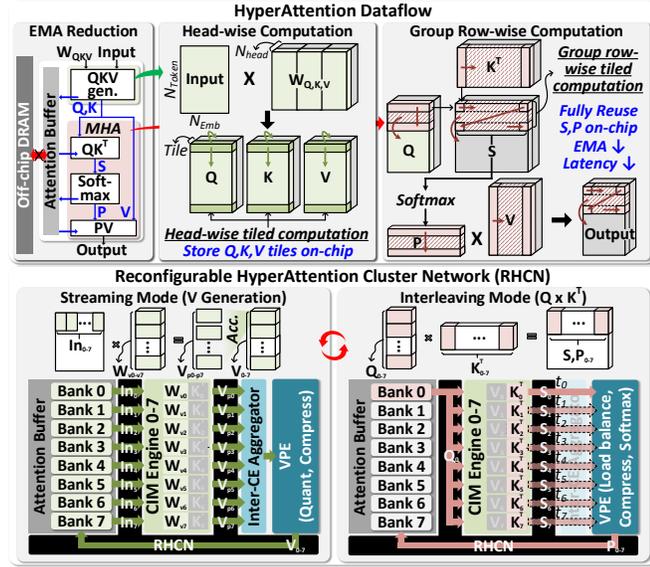


Fig. 3. HyperAttention (HA) dataflow and RHCN to reduce EMA

optimal data movements for HA. RHCN reconfigures the on-chip network for streaming mode and interleaving mode. The streaming mode enables high BW data fetch from the attention buffer to VPE through CEs and ICEA, while the interleaving mode broadcasts data from a single bank of the attention buffer to CEs and interleaves CEs to fetch S to VPE for softmax. The IMP, enabled by the area-efficient dual-port 3T cells, dual X-DEC, and MAC latch, concurrently processes read, write, and MAC, reducing W fetching and refresh latency by 66.6% and 49.9%, respectively, as shown in Fig. 4. This concurrent read/write/MAC is achieved by only 3.33 transistors per cell, 1.9 \times lower than [5], [6]. Thanks to RHCN and IMP, HA ultimately allows off-chip access only for input and output of attention, achieving up to 6.2 \times EMA and 9.8 \times latency

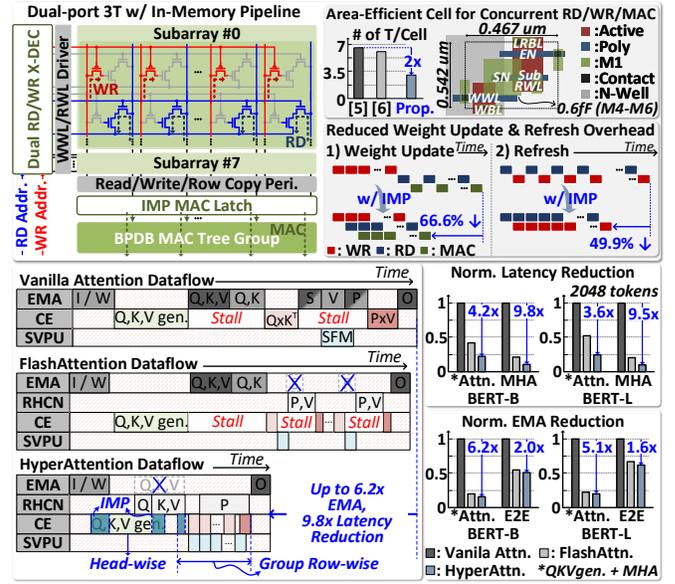


Fig. 4. Dual-port 3T with in-memory pipeline (IMP) and HA dataflow

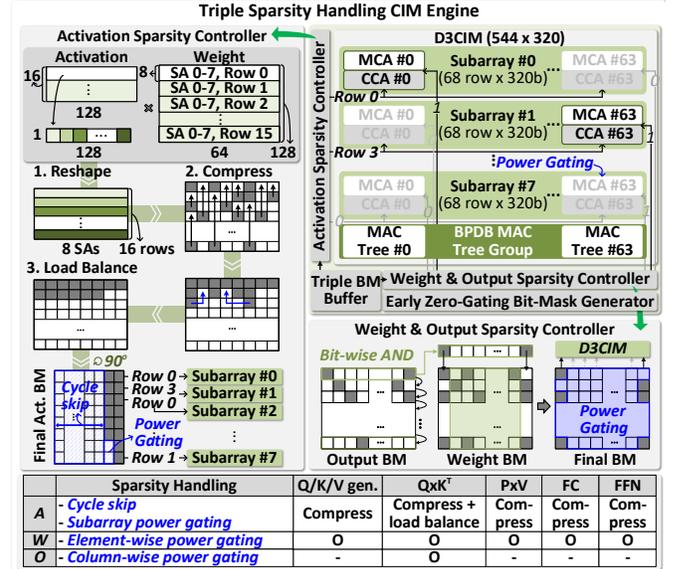


Fig. 5. Triple sparsity handling CIM engine for activation, weight, and output

reduction over vanilla attention.

B. Triple Sparsity Handling CIM Engine

Fig. 5 demonstrates the TSH for A/W/O sparsity to increase computation efficiency. CE handles sparsity in activations for cycle skip and subarray gating, weights for element-wise gating, and output for column-wise gating. Based on the activation BM, pre-processed in the compressor, the TSC only activates nonzero subarrays of the D3CIM to reduce energy and separately controls the row address of each subarray to skip computation cycles. To skip cycles further, the TSC applies load balancing by dynamically redistributing the W with the row copy. The copied W can be reused for subsequent computations based on the column-wise spatial locality in

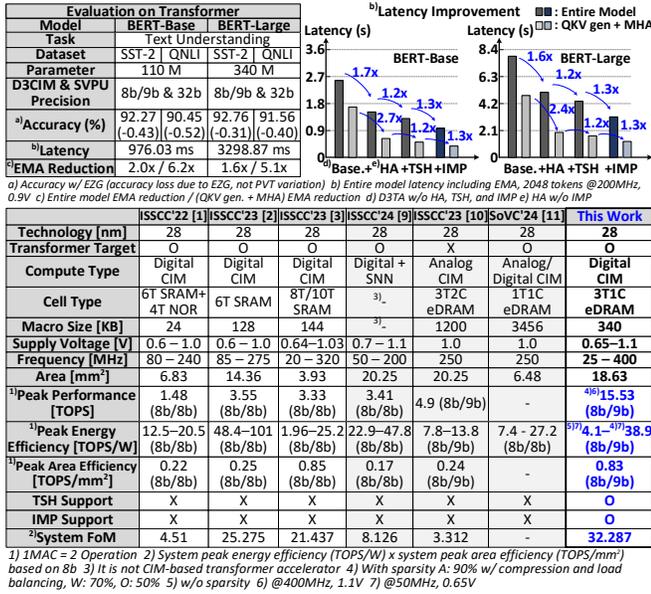


Fig. 10. Measurement results and performance comparison table

48b encoded data broadcasts to 64 columns and is multiplied with 40b W data through the 1b-shifter-based booth multiplier in each column. If ONE and TWO are zero, multiplication is skipped. The final outputs are 64 11b, considering the bit position. The BPDB MAC tree group achieves 3.86 \times and 2.38 \times higher area-efficiency compared to bit-series and bit-parallel computing.

III. MEASUREMENT RESULTS AND CONCLUSION

The D3TA operates at 25-to-400MHz with 0.65-to-1.1V, and by leveraging key features, it achieves up to 2.6 \times latency reduction and 2.0 \times EMA reduction compared to the baseline, as shown in Fig. 10. The system peak throughput is 15.53 TOPS with energy efficiency of 38.9 TOPS/W, and the system FoM is 1.28-to-9.75 \times higher than [1]–[3], [9], [10]. Fig. 11 shows the chip photograph and performance summary of D3TA fabricated in 28nm CMOS technology. In conclusion, D3TA optimizes transformer acceleration in memory, achieving high energy and area efficiency with reduced EMA and latency.

ACKNOWLEDGMENT

This work was supported by IITP grant funded by the Korea government (MSIT) (No.2022-0-01037, Development of High Performance Processing In Memory Technology based on DRAM and IITP-2025-RS-2023-00256472, Graduate School of Artificial Intelligence Semiconductor).

REFERENCES

- [1] F. Tu *et al.* "A 28nm 15.59 μ J/token full-digital bitline-transpose CIM-based sparse transformer accelerator with pipeline/parallel reconfigurable modes" 2022 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2022.
- [2] F. Tu *et al.* "16.1 MuITCIM: A 28nm 2.24 μ J/token Attention-Token-Bit Hybrid Sparse Digital CIM-Based Accelerator for Multimodal Transformers." 2023 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2023.

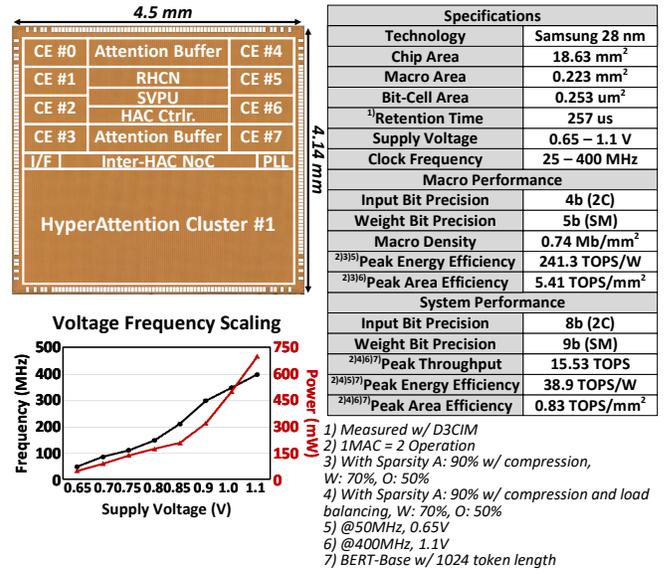


Fig. 11. Chip photograph and performance summary

- [3] S. Liu *et al.* "16.2 A 28nm 53.8 TOPS/W 8b sparse transformer accelerator with in-memory butterfly zero skipper for unstructured-pruned NN and CIM-based local-attention-reusable engine" 2023 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2023.
- [4] D. Kim *et al.* "DPIM: A 19.36 TOPS/W 2T1C eDRAM Transformer-in-Memory Chip with Sparsity-Aware Quantization and Heterogeneous Dense-Sparse Core" 2024 European Solid-State Electronics Research Conference (ESSERC). IEEE, 2024.
- [5] W.-S. Khwa *et al.* "34.2 A 16nm 96Kb Integer/Floating-Point Dual-Mode-Gain-Cell-Computing-in-Memory Macro Achieving 73.3-163.3 TOPS/W and 33.2-91.2 TFLOPS/W for AI-Edge Devices" 2024 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2024.
- [6] H. Fujiwara *et al.* "34.4 A 3nm, 32.5 TOPS/W, 55.0 TOPS/mm² and 3.78 Mb/mm² Fully-Digital Compute-in-Memory Macro Supporting INT12 \times INT12 with a Parallel-MAC Architecture and Foundry 6T-SRAM Bit Cell" 2024 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2024.
- [7] A. Yazdanbakhsh *et al.* "Sparse attention acceleration with synergistic in-memory pruning and on-chip recomputation" 2022 55th IEEE/ACM International Symposium on Microarchitecture (MICRO). IEEE, 2022.
- [8] Y.-D. Chih *et al.* "16.4 An 89TOPS/W and 16.3 TOPS/mm² all-digital SRAM-based full-precision compute-in memory macro in 22nm for machine-learning edge applications" 2021 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2021.
- [9] S. Kim *et al.* "20.5 C-Transformer: A 2.6-18.1 μ J/Token Homogeneous DNN-Transformer/Spiking-Transformer Processor with Big-Little Network and Implicit Weight Generation for Large Language Models" 2024 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2024.
- [10] S. Kim *et al.* "16.5 DynaPlasia: An eDRAM in-memory-computing-based reconfigurable spatial accelerator with triple-mode cell for dynamic resource switching" 2023 IEEE International Solid-State Circuits Conference (ISSCC). IEEE, 2023.
- [11] S. Hong *et al.* "Dyiamond: A 1T1C DRAM In-memory Computing Accelerator with Compact MAC-SIMD and Adaptive Column Addition Dataflow" 2024 IEEE Symposium on VLSI Technology and Circuits (VLSI Technology and Circuits). IEEE, 2024.