

A 28nm 244.45TOPS/W Winograd-Standard Fusion accelerator with Symmetric Hybrid Domain CIM Groups for Edge AI devices

An Guo^{1#}, Zhichao Liu^{1#}, Wentao Zheng¹, Yutong Zhang¹, Tianhui Jiao¹, Fangyuan Dong¹, Mingzi Wang², Shaochen Li¹, Zhican Zhang¹, Yuhui Shi¹, Xing Wang¹, Xin Si¹, Xin Wang^{2*}, Wenwu Zhu^{2*}

¹Southeast University, Nanjing, China, ²Tsinghua University, Beijing, China

Email: xin_wang@tsinghua.edu.cn/, wwzhu@tsinghua.edu.cn/

#Equally contributed authors

*Corresponding authors

The rapid advancement of AI technologies has garnered widespread attention, but the substantial power consumption of sophisticated AI models presents a significant challenge. Various AI accelerators [1–10] have been developed to optimize energy efficiency for AI tasks. Winograd Convolution, as a CNN accelerating method, can save more than 2.25 times operations, as illustrated in Fig. 1. This method involves four steps: (1) Winograd-domain (WD) weights transformation to W' using the G matrix; (2) WD features transformation to F' using the B matrix; (3) Hadamard production between W' and F' to obtain WD results Y' ; and (4) Inverse WD transformation into final neural network results Y . Winograd Convolution typically employs multiple transform sizes, such as (2,3), (4,3), (6,3), where the numbers represent output and input tile sizes (OTS and ITS) respectively. Larger OTS yields more operation savings but increases accuracy loss, with OTS >4 leading to task failure in YOLO inference tests. Given the additional memory requirements of Winograd convolution, CIM emerges as a promising solution to efficiently implement this algorithm. This work proposes a Winograd AI accelerator with symmetric digital-analog hybrid domain CIM groups addressing: (1) Winograd's limitation to stride-1 layers, as decomposition algorithms for stride-2 CONV lead to large redundant calculations; (2) Significant storage and access overhead in WD weights; and (3) Complex Winograd operations, including GWGT, BTFB, Hadamard production and ATY'A. A 28nm fabricated Winograd/standard fusion AI accelerator achieves an INT8 energy efficiency of 79.9TOPS/W and a mean average precision (mAP) of 57.88% with retraining for YOLOv8 with Winograd CONV.

Fig. 2 illustrates the architecture of the proposed Winograd and standard convolution fusion (WSF) hybrid CIM chip, featuring: (1) a WSF CNN CIM architecture for edge image and video processing; (2) a Left Feature Parallel, Right Feature Serial Adder (LFP-RFsA) and WSF-digital-to-macro buffer (WSF-D2MB) based pre-macro processor for input stationary flow; (3) a bit-configurable digital-analog hybrid domain CIM with symmetric structure (S-hybrid CIM); and (4) a Winograd-advanced-calculation standard-ordinary-accumulation (WaC-SoA) based post-macro processor. The WSF chip comprises a main control, IO, Input preprocessor, 568Kb global buffer, SIMD core for nonlinear functions, WSF mode controller, ReLU early termination, and WSF core. Winograd calculations are divided into four parts: (1) $F' = BTFB$, computed in the LFP-RFsA pre-macro processor, where input data is split into LFP- and RFs-groups, processed in parallel and bit-serial over 128 cycles respectively, with WSF-D2MB using paired DFFs for pipeline input scheduling; (2) $W' = GTWG$, calculated offline for multi-time reuse; (3) $Y' = F' \odot W'$, executed in 4 symmetric-hybrid CIMs capable of INT2/4/8 mode operations, solving 128-ichs MAC operations; and (4) $Y = ATY'A$, processed in the WaC-SoA aggregator, computing ATY' and half of Y'A in advance to accelerate post-macro processing. Final results are sent to an output buffer for subsequent ReLU and other nonlinear functions.

Fig. 3 illustrates the proposed LFP-RFsA and WSF-D2MB pre-macro processor, which operates in two modes: Winograd and standard convolution. In Winograd mode, input data undergoes BTFB transformation, allowing multi-position input sharing. The process is divided into two parts, calculated alternately. Each part shares 4 of 12 inputs between left- and right-groups, with different input rows directed to 4 S-hybrid CIM macros. The LFP-RFs scheduler manages left-groups using 2 8b individual and 2 8b shared inputs for one LFP-adder per macro, while right-groups use 8x2 1b individual and 8x2 1b shared inputs for 8 RFs-adders per macro. LFP-adders perform 4

8b accumulations or subtractions per cycle, producing 8b macro input data stored in a temporary storage (TS) buffer. RFs-adders calculate 4 1b operations per cycle, with 8 adders introduced to maximize utilization of CIM macros' 1b feature and 8b weight MAC capabilities for 16-ochs in 16 cycles. RFs-adder results are stored in TS buffers with 3b carry bit storage DFF and 1b extra DFF for the last 1b PSUM. To address carry-in issues, RFsAs operate two cycles earlier in the Winograd BTFB pipeline pre-macro flow. In standard convolution mode, inputs from different ochs are divided into left- and right-groups, with left features scanned from high to low bits and right features from low to high. These inputs bypass LFP-RFsA and WSF-D2MB, proceeding directly to S-hybrid CIM macros.

Fig. 4 illustrates the proposed symmetric hybrid CIM macro, comprising an INDRV, WL driver, key switch driver, IO, multi-level accumulator, main and CIM controllers, and 8 128x32 banks. Each bank processes 128 1b and 8b MAC operations for 2 output channels (ochs), one in analog and one in digital domain. The macro employs an input-stationary inner-macro data flow, maximizing input utilization through local weight switching. The left side of the figure details the computing units, where upcells and downcells connect to digital computing units (DCU) and analog computing units (ACU), controlled by K0. Different banks are managed by distinct Kns, as shown in the top right. The 8b vertical cut example operates in 8 steps: (1) With $K0-6=0$ and $K7=1$, left features input highest bits and right features input lowest bits. Left features and upcells compute digitally, except for the lowest bit in analog, while right features and downcells compute in analog, except for the highest bit in digital. This process repeats for 16 cycles, switching local weights for 16 ochs. (2) $K0-5=0$ and $K6-7=1$, inputting MSB-1 for left features and LSB+1 for right features. Left features and upcells compute digitally except for the lowest 2 bits in analog, and right features and downcells compute in analog except for the highest two bits in digital. (3) Similar operations as (1-2) continue with Kn and left-right feature switches, completing the 8-step process for full 8b computation.

Fig.5 top illustrates the proposed WaC-SoA post-macro processor, utilizing S-hybrid CIM macro in INT4/8 mode. The S-hybrid CIM macro incorporates two multi-level accumulators, generating two 128 1b feature and 4b weight PSUMs, which are processed by distinct Winograd advanced ATY'A processors. Each WaC processor divides calculations into top- and bottom- A matrix multiplications, computing ATY' & top of Y'A initially and ATY' & bottom of Y'A subsequently, with shared input data. In INT8 mode, post-WaC processor PSUMs undergo shift-and-add operations to produce final INT8 results. Fig.5 bottom showcases performance improvements: the LFP-RFs pre-macro processor achieves 64x adders, 1.79x power, and 4.59x area savings, while the WaC-SoA Aggregator post-macro processor realizes 1.50x power reduction and 3.04x speed increase. The CIM macro achieves an area efficiency of 840.9GOPS/mm²@0.9V and an energy efficiency of 79.88TOPS/W@0.6V. Three hybrid domain CIM structures - lightning-like (LLS) [5], vertical-cut (VCS) [6], and symmetric - were evaluated using 1000 random vectors with identical digital and analog circuits, resulting in 2-8x RC overhead with maximum 4.01% accuracy loss. In the WSF chip, pre-, inner-, and post-macro components consume 7.3%, 20.5%, and 72.2% of power and 10.7%, 17.6%, and 71.7% of area, respectively, while conducting 4.2%, 8.2%, and 87.6% of operations. The additional operations in pre- and post-macros are attributed to BTFB and ATY'A computations.

Fig.6 displays the measured shmoo plot of the proposed WSF-chip, along with measured results on VGG16, ResNet18, and Yolov8, and a comparison table with previous work. The chip achieves 620MHz@0.9V and 165MHz@0.6V, with an energy efficiency of 45.46-243.62TOPS/W. Fig.7 presents the die photo and chip summary of the fabricated 28nm CMOS technology Winograd-standard CNN fusion Hybrid CIM accelerator. Proposed chip features a core area of 0.9mm² with four 0.11mm² hybrid CIM macros, achieving an energy efficiency of 62.93-244.45 TOPS/W @INT4 precision and 14.85-57.68 TOPS/W @INT8 precision. Additionally, it demonstrates an inference accuracy of 70.72% on ResNet-18@ImageNet and a mean average precision (mAP) of 57.88% on Yolov8@COCO.

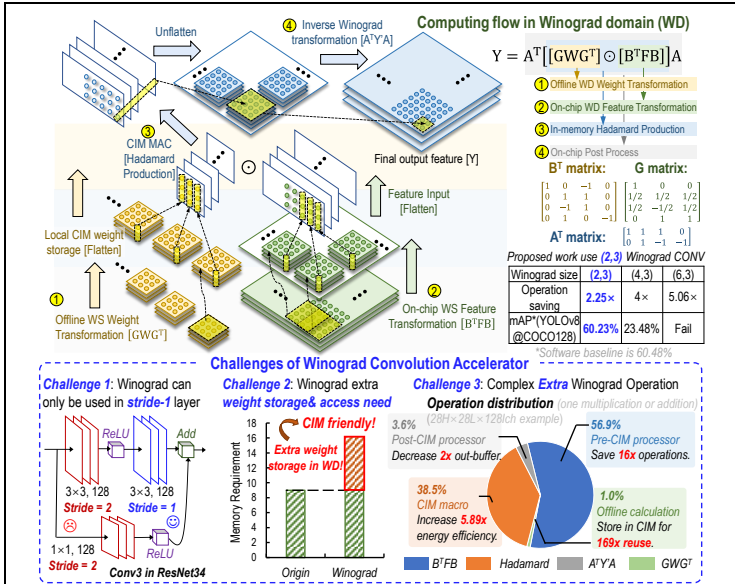


Fig. 1. Design challenges of Winograd CIM accelerator.

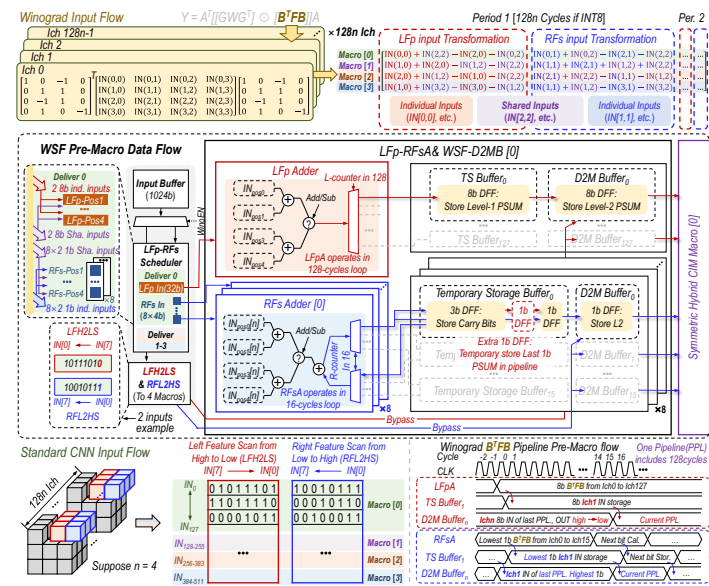


Fig. 3. Proposed LFP-RFSA and WSF-D2MB pre-macro processor.

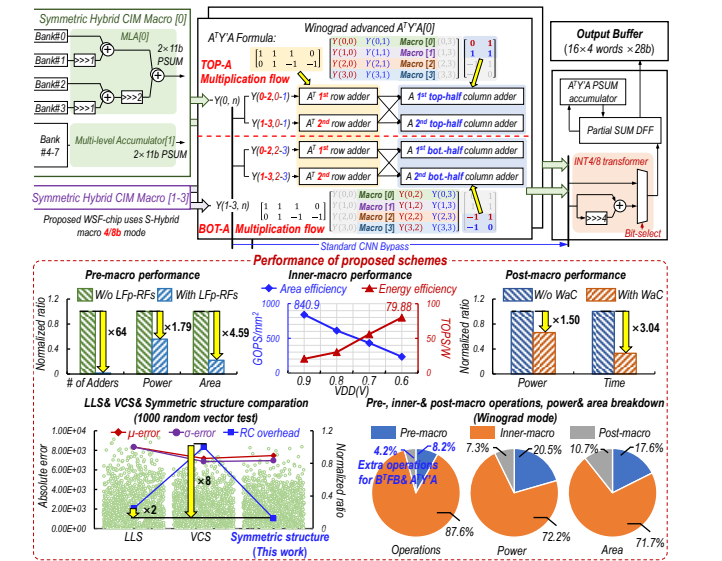


Fig. 5. Proposed WaC-SoA post-macro processor and performance improvements.

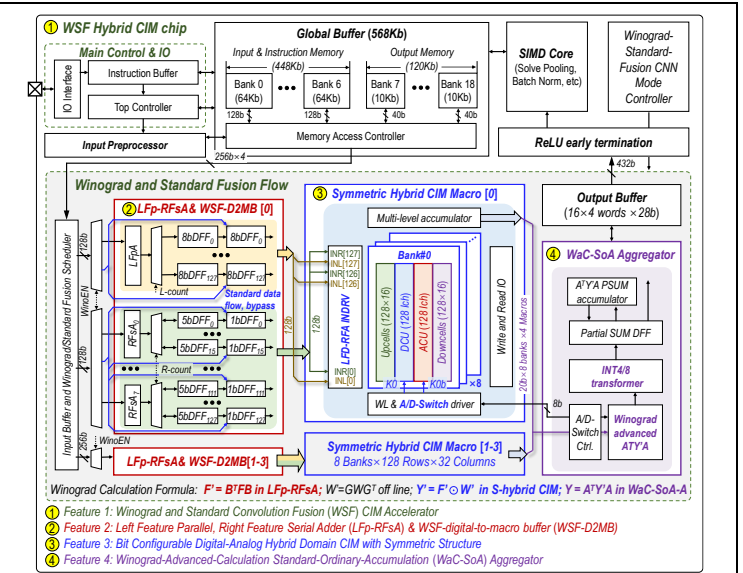


Fig. 2. Overall system architecture of proposed Winograd-Standard convolution fusion CIM accelerator.

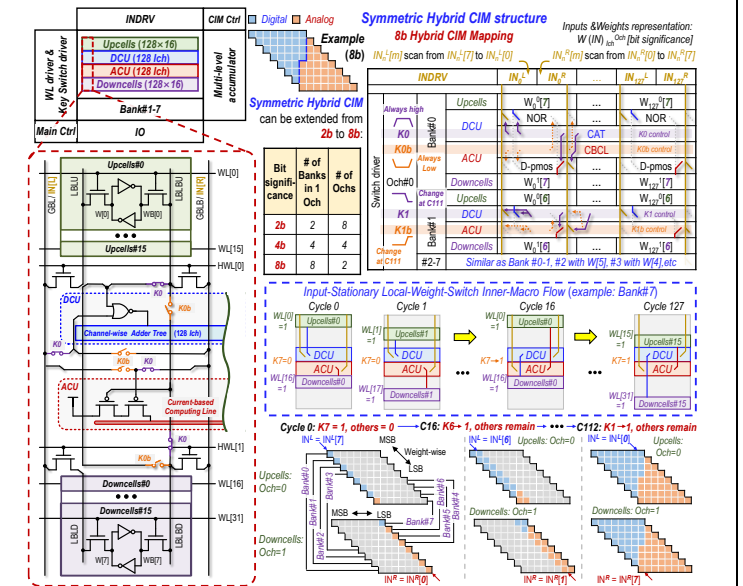


Fig. 4. Proposed symmetric hybrid CIM macro.

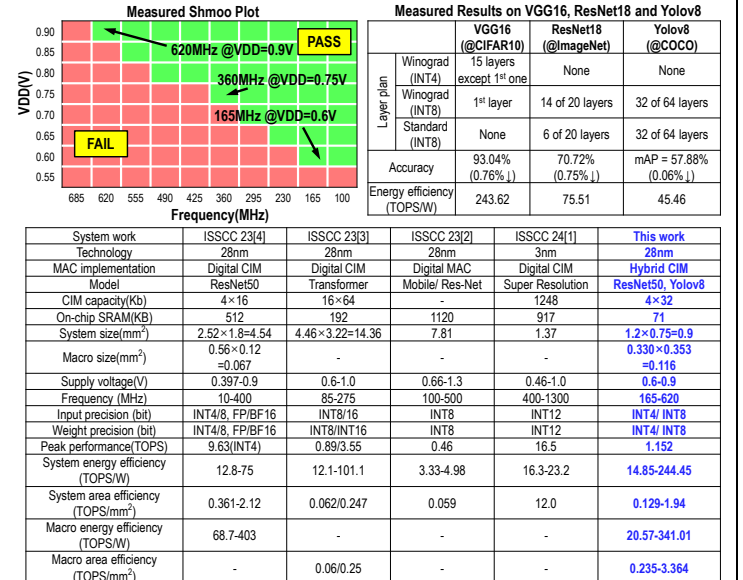
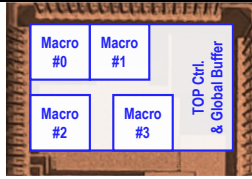


Fig. 6. Measurement results and performance comparison table.



¹High point: test using YOLOv8@COCO under 0.6V at 165MHz; low point: test using YOLOv8@COCO under 0.9V at 650MHz.

²High point: test under 0.9V; low point: test under 0.6V.

³Using ResNet-18 model and the software baseline is 71.27%.

⁴Using YOLOv8 model and the software baseline is 57.94%.

CHIP SUMMARY		
Technology	28nm CMOS	
Core area	1.2×0.75=0.9 mm ²	
CIM Macro Area	4×0.330×0.353 mm ²	
Supply voltage(V)	0.6-0.9	
Frequency	165-620MHz	
CIM size	4×32Kb	
SRAM size	568Kb	
CIM memory density(Kb/mm ²)	90.65	
Peak performance(TOPS)	1.152	
Precision(bit)	INT 4	INT 8
System energy efficiency ¹ (TOPS/W)	62.93-244.45	14.85-57.68
Macro energy efficiency ¹ (TOPS/W)	93.18-341.01	20.57-79.9
System area efficiency ² (TOPS/mm ²)	0.516-1.94	0.129-0.485
Macro area efficiency ² (TOPS/mm ²)	0.94-3.364	0.235-0.841
Neural network index of different AI tasks	Inference accuracy ³ (ResNet-18 @ImageNet)	70.72% (16, 4 of total 20 layers using Winograd INT8 and Standard INT8 convolution)
	mAP ⁴ (YOLOv8@COCO)	57.88% (32, 32 of total 64 CONV layers using Winograd INT8 and Standard INT8 convolution)

Acknowledge

This work is funded with NSTMP under Grant 2022ZD0118901, NSFC under Grant 62522403, 62204036, 92464302, Jiangsu Provincial RDP under Grant BE20230201 and the Fundamental Research Funds for the Central Universities.

References:

[1] M. Shih et al., "NVE: A 3nm 23.2TOPS/W 12b-Digital-CIM-Based Neural Engine for High-Resolution Visual-Quality Enhancement on Smart Devices," ISSCC, pp. 128-130, 2024.

[2] C. Du et al., "A 28nm 11.2TOPS/W Hardware-Utilization-Aware Neural-Network Accelerator with Dynamic Dataflow," ISSCC, pp. 332-333, 2023.

[3] F. Tu et al., "MuTCIM: A 28nm 2.24μJ/Token Attention-Token-Bit Hybrid Sparse Digital CIM-Based Accelerator for Multimodal Transformers," ISSCC, pp. 248-249, 2023.

[4] J. Yue et al., "A 28nm 16.9-300TOPS/W Computing-in-Memory Processor Supporting Floating-Point NN Inference/Training with Intensive-CIM Sparse-Digital Architecture," ISSCC, pp. 252-254, 2023.

[5] A. Guo et al., "A 22nm 64kb Lightning-Like Hybrid Computing-

Fig. 7. Die photo and chip summary table.

in-Memory Macro with a Compressed Adder Tree and Analog-Storage Quantizers for Transformer and CNNs," ISSCC, pp. 570-571, 2024.

[6] P. -C. Wu et al., "A 22nm 832Kb Hybrid-Domain Floating-Point SRAM In-Memory-Compute Macro with 16.2-70.2TFLOPS/W for High-Accuracy AI-Edge Devices," ISSCC, pp. 126-128, 2023.

[7] A. Guo et al., "A 28nm 64-kb 31.6-TFLOPS/W Digital-Domain Floating-Point-Computing-Unit and Double-Bit 6T-SRAM Computing-in-Memory Macro for Floating-Point CNNs," JSSC, vol.59, no.9, 2024.

[8] V. Chikin and V. Kryzhanovskiy, "Channel Balancing for Accurate Quantization of Winograd Convolutions," CVPR, pp. 12497-12506, 2022.

[9] C. Mu et al., "A 28nm 76.25TOPS/W RRAM/SRAM-Collaborative CIM Fine-Tuning Accelerator with RRAM-Endurance/Latency-Aware Weight Allocation for CNN and Transformer," 2024 IEEE Asian Solid-State Circuits Conference (A-SSCC), Hiroshima, Japan, 2024, pp. 1-3.

[10] Z. Wei, Y. Su, T. T. -H. Kim and Y. Zheng, "A 280.580 μm²/Grid, 30.7pJ TMRA-PIM (Time-Multiplexed Random-Access Processing-in-Memory) Programmable Accelerator for Partial Differential Equations Solving and Edge Processing," 2024 IEEE Asian Solid-State Circuits Conference (A-SSCC), Hiroshima, Japan, 2024, pp. 1-3.