TRIFP-DCIM: A Toggle-Rate-Immune Floatingpoint Digital Compute-in-Memory Design with Adaptive-Asymmetric Compute-Tree

Xing Wang, Tianhui Jiao, Shaochen Li, Yuchen Ma, Zhican Zhang, Zhichao Liu, Xi Chen and Xin Si* ASP-DAC 2025

January 23, 2025

Challenges and proposed solutions

Current FP-DCIM suffers from:

- 😥 Large area overhead of mantissa alignment
- Large-scale adder tree limits the energy/area efficiency
- Less sparsity utilization

Solutions in this work

TRIFP-DCIM proposes:

- Serial Shift Scheme for BF16 MAC
- Adaptive Asymmetric Compute-tree(AACT) Circuit
- TrifectaOne method of increasing data sparsity

TRIFP-CIM Overall Architecture



Mantissa needs to be processed and expanded to 9 bits, adding a hidden "1" and a sign bit.

- Mantissa MAC Block (MMACB)
- Shift and Accumulate Block (SAB)
- Exponent Accumulation Block (EAB)
- Mantissa Shift Block (MSB)

BF16 Computation Dataflow



- 1) Compute 9-bit exponent sums from 8-bit exponents.
- 2) Find maximum exponent using a comparator tree.
- 3) Serial shift mantissa based on control signals.
- 4) Perform multiply-accumulate to get PMACV.
- 5) Integrate PMACV and exponent into BF16 format.

Serial Shift Module

9+2 bits Shift Register Stack



- Clock Gating Units & Shift Register Stack (SRS)
- A SR is made up of 11 serially connected registers
- Highest 2 bits sign bit extension, lower 9 bits mantissa
- Deliver 2b of the mantissa to subsequent computation circuit each cycle



- 1) Given *E1=011101010* and *E2=011100111*.
- 2) The counter counts down from 01110101.
- 3) The counter value matches *E1* at *T0*, then first SR starts shifting.

When the counter value aligns with E1, the first SR enters the "Find" state at T0.



- 1) Given *E1=011101010* and *E2=011100111*.
- 2) The counter counts down from 01110101.
- 3) The counter value matches *E1* at *T0,* then first SR starts shifting.
- 4) The counter value continues to decrease by one at *T1* and *T2*.

Consequently, the first SR starts shifting at T1, while the second SR remains in the "Stay" state.



- 1) Given *E1=011101010* and *E2=011100111*.
- 2) The counter counts down from 01110101.
- 3) The counter value matches *E1* at *T0,* then first SR starts shifting.
- 4) The counter value continues to decrease by one at *T1* and *T2*.

The second SR transitions to the "Find" state at T2 when the counter value equals E2, and the first SR remains in the "Shift" state.



- 1) Given *E1=011101010* and *E2=011100111*.
- 2) The counter counts down from 01110101.
- 3) The counter value
- matches *E1* at *T0,* then first SR starts shifting.
- 4) The counter value continues to decrease by one at *T1* and *T2*.
- 5) The second SR starts shifting at *T3.*
- 6) Continue shifting until the specified cycle number is reached.

Toggle Path Research



Difference:

Position of two toggling inputs

Situation 1: Close Toggles share a data path cross at the first adder level

Situation 2: Far Away Toggles propagate along respective paths cross at the fourth adder level

The power consumption of situation1 is about 50% less than that of situation2 under certain input patterns.

Toggle Path Research



Conclusion: gathering toggling inputs to one side can lead to reduction of power consumption. 11

TrifectaOne Method



- Concentrate the toggles on one side of the compute-tree.
- Ensure 4-2 compressors on the other side will not affect the inference accuracy.

Weight Asymmetric Sparsification

Algorithm 1 Multi-channel weight asymmetric sparsification.

Input :Four-dimensional tensor *W*₀ **Output**:Four-dimensional tensor *W*_a

- // Split channel
- 1 **if** *input_channel < 8* **then**
- 2 $W_a = W_0;$
- 3 end
- 4 **if** *input_channel* > 8 **then**
- 5 split $W_0 \rightarrow$ chunks[N-1:0], N = Quotient of *input_channel/n_size*;

6 end

7 **for** *i* = 0 to (N-1) **do** if i%2 = 1 then // all data in chunks[i] 8 $S_c = \text{Sign}(chunks[i]), A_c = \text{Abs}(chunks[i]);$ 9 // Mantissa 3b Approximation $exp0 = \lfloor \log_2(A_c) \rfloor, temp = 2^{exp0},$ 10 $R_c = A_c - temp;$ $exp1 = \lfloor \log_2(R_c) \rfloor, temp = temp + 2^{exp1},$ 11 $R_c = A_c - temp;$ $exp2 = \lfloor \log_2(R_c) \rfloor, temp = temp + 2^{exp2};$ 12 $chunk[i] = temp \times S_c;$ 13 end 14 15 **end**

16 Merge chunks on dim = 1, W_a equal to merged result; 17 return W_a

Step1:

Group 8 input channels as a chunk, alternately set as sparse chunks and dense chunks.

Step2:

Increase the sparsity of data in the sparse chunk by retaining only the first 3 ones of the mantissa.

Adaptive-asymmetric Compute-Tree



16 2-bit numbers input

Toggle-intensive side: Full-precision CT

Toggle-sparse side:

Approximate CT

- 4-2 compressors
- approximate 2bitadders

Asymmetric weight sparsity Adaptive Asymmetric Compute-Tree Circuit

Adaptive-asymmetric Compute-Tree



Adaptive-asymmetric Compute-Tree



Test result of TRIFP-DCIM



Accuracy test under different networks

Sparsity differences test



Power consumption test of the TrifectaOne alone with AACT

Comparison With Previous Works

	JSSCC		DAC'21[2]	ISSCC ²³ [3]	ISSC	CC 23[4]	This Work
A most size	512×128		256×64	02017h	2×512×128		$144 \times 512 + 648 \times 128$
Allay size	(64Kb)		(16Kb)	032KD	(128Kb)		(153Kb)
Technology	28nm		28nm	22nm	28nm		28nm
Cell Structure	6T+LCC		8T+Local Peri	6T+DAH-MCA	DBcell+FCU		6T+AACT
MAC Operation	Analog		Analog	Hybrid	Digital		Digital
Supply Voltage	0.7-0.9		0.5-1.1	0.6-0.8	0.6-0.9		0.6-0.8
Precision(bit)	4	8	¹ 2/8	BF16	8	BF16	BF16
Energy Efficiency (TOPS/W)	47.8	11.5	8.78-36.10	16.22-17.59 TFLOPS/W	19.5-44	14.04-31.6 TFLOPS/W	² 14.51- ³ 36.83 TFLOPS/W

¹ The input precision is 2bit, and the weight precision is 8 bit.

² Simulation results with FP worst case under 0.8V.

³ Simulation results with actual ResNet18 model @CIFAR100 under 0.6V.

Combined with the **TrifectaOne method** and **AACT design**, the proposed FP-CIM is simulated under 28nm technology process.

Compared with other works, this work achieves a high energy efficiency of **14.51-36.83 TFLOPS/W** in the digital domain.

Conclusion

TRIFP-DCIM:

28nm digital-domain CIM with **BF16 precision**;



Achieve 14.51-36.83TFLOPS/W@BF16.

Reference

[1] Xin Si et al. A local computing cell and 6t sram-based computing-in memory macro with 8-b mac operation for edge ai chips. IEEE Journal of Solid-State Circuits, 56(9):2817–2831, 2021.

[2] Kyeongho Lee et al. A charge-sharing based 8t sram in-memory computing for edge dnn acceleration. In 2021 58th ACM/IEEE Design Automation Conference (DAC), pages 739–744, 2021.

[3] Ping-Chun 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. In 2023 IEEE International Solid-State Circuits Conference (ISSCC), pages 126–128, 2023.

[4] An Guo et al. A 28nm 64-kb 31.6-tflops/w digital-domain floating pointcomputing-unit and double-bit 6t-sram computing-in-memory macro for floating-point cnns. In 2023 IEEE International Solid-State Circuits Conference (ISSCC), pages 128–130, 2023.

THANK YOU