In-Memory Data Parallel Processor (ASPLOS'18)

Author: Daichi Fujiki(<u>dfujiki@umich.edu</u>) Presenter: Choi, Jinyoung

PIM (Processing in Memory)

- Memory access is the main bottleneck in modern Von-Neumann • architecture. 1000
 - CPU-memory performance gap
 - Off-chip access is expensive!



 People think about implanting computing logic into the memory unit.





PRIME (ISCA'16)





Figure 3. (a) Traditional shared memory based processor-coprocessor architecture, (b) PIM approach using 3D integration technologies, (c) PRIME design.





RSITY OF CALIFORNIA



а

Dot operation using Re-RAM

$$b = \begin{bmatrix} a_1 & a_2 & a_3 & a_4 & a_5 \end{bmatrix}_{(1 \times n)} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{bmatrix}_{(n \times 1)} = \left\{ a_1 b_1 + a_2 b_2 + a_3 b_3 + a_4 b_4 + a_5 b_5 \right\}$$

$$V1 = \left\{ a_1 b_1 + a_2 b_2 + a_3 b_3 + a_4 b_4 + a_5 b_5 \right\}$$

$$V2 = Dot Product$$

Basic Laws of Circuits







(a) Multiply-Accumulate operation

(b) Vector-Matrix Multiplier

Fig. 1. (a) Using a bitline to perform an analog sum of products operation. (b) A memristor crossbar used as a vector-matrix multiplier.

(ISAAC, ISCA'16)

UC REVERSITY OF CALIFORNIA



Why this paper?

- In previous Re-RAM accelerators, •
 - they support only dot-product operation.
 - \cdot we need to manually map specialized kernel to the memory arrays.

latency. In this paper, we seek to answer the question, to what extent is resistive memory useful for more general-purpose computation?





Proposed architecture: In-Memory Data Parallel Processor

- In previous Re-RAM accelerators, •
 - they support only dot-product operation. •
 - ISA beyond dot-product •
 - we need to manually map specialized kernel to the memory arrays.
 - Compiler •





Architecture & ISA





Figure 1. In-Memory Processor Architecture. (a) Hierarchical Tiled Structure (b) ReRAM array Structure



- XB
- ADC
- S/A

Opcode	Format	Cycles
add	<mask><dst></dst></mask>	3
dot	<mask><reg_mask><dst></dst></reg_mask></mask>	18
mul	<src><src><dst></dst></src></src>	18
sub	<mask><mask><dst></dst></mask></mask>	3
shift{l r}	<src><dst><imm></imm></dst></src>	3
mask	<src><dst><imm></imm></dst></src>	3
mov	<src><dst></dst></src>	3
movs	<src><dst><mask></mask></dst></src>	3
movi	<dst><imm></imm></dst>	1
movg	<gaddr><gaddr></gaddr></gaddr>	Variable
lut	<src><dst></dst></src>	4
reduce_sum	<src><gaddr></gaddr></src>	Variable





UNIVERSITY OF CALIFORNIA









Figure 6. IB Expansion.

Execution model



Figure 3. Execution Model.





Input matrix A

Input matrix B

Data Flow Graphs

Modules

ReRAM arrays

UNIVERSITY OF CALIFORNIA

Testbed

- Baseline 1: Intel Xeon E5-2697 v3 multi-socket server •
- Baseline 2: Nvidia Titan XP
- In-Memory Processor(IMP)

Parameter	CPU (2-sockets)	GPU (1-card)	IMP
SIMD slots	448	3840	2097152
Frequency	3.6 GHz	1.58 GHz	20 MHz
Area	912.24 mm^2	$471 \ mm^2$	$494 \ mm^2$
TDP	290 W	250 W	416 W
Memory	7MB L2; 70MB L3	3MB L2	1GB
	64GB DRAM	12GB DRAM	RRAM

Table 5. Comparison of CPU, GPU, and IMP Parameters



Component	Params	Spec	Power	Area(mm ²)
ADC	resolution	5 bits	64 mW	0.0753
	frequency	1.2 GSps		
	number	64×2		
DAC	resolution	2 bits	0.82 mW	0.0026
	number	64×256		
S+H	number	64×128	0.16 mW	0.00025
ReRAM	number	64	19.2 mW	0.0016
Array				
S+A	number	64	1.4 mW	0.0015
IR	size	2KB	1.09 mW	0.0016
OR	size	2KB	1.09 mW	0.0016
Register	size	3KB	1.63 mW	0.0024
XB	bus width	16B	1.51 mW	0.0105
	size	10×10		
LUT	number	8	6.8 mW	0.0056
Inst. Buf	size	8×2 KB	5.83 mW	0.0129
Router	flit size	16	0.82 mW	0.00434
	num_port	9		
S+A	number	1	0.05 mW	0.000004
I Tile Total			101 mW	0.12
Inter-Tile	number	584	0.81 W	2.50
Routers				
Chip total			416 W	494 mm ²

Table 4. In-Memory Processor Parameters

Operations







Figure 9. Multiplication Latency.









before, in-memory accelerator is better coupled with the existing memory hierarchy to avoid data loading overhead.



□ Kernel ■ Data loading ■ NoC ■ Sequential+Barrier CPU IMP IMP CPU IMP IMP CPU IMP IMP Mem Mem Mem Fluidanimate Streamcluster Canneal

Figure 12. CPU Application performance.



Power & Energy







Compiler optimization



MaxDLP: one IB for a module - cannot exploit ILP - under utilization with small small size MaxILP: maximum IBs for a module - serialization overhead MaxArrayUtil: considering #of SIMD slots



- need to invoke kernel multiple time (SIMD slot < data)
 - needed by input data



Thanks

• QnA



UC RIVERSITY OF CALIFORNIA



Re-RAM array operations







ation (d) Subtraction

UC REPORT OF CALIFORNIA