

# Newton: A DRAM-maker's Accelerator-in-Memory (AiM) Architecture for Machine Learning

MICRO'20

Mingxuan He  
*Electrical and Computer Engineering*  
*Purdue University*  
West Lafayette, IN, U.S.A.  
he238@purdue.edu

Choungki Song  
*DRAM Design*  
*SK Hynix*  
Icheon, South Korea  
choungki.song@sk.com

Ilkon Kim  
*DRAM Design*  
*SK Hynix*  
Icheon, South Korea  
ilkon.kim@sk.com

Chunseok Jeong  
*DRAM Design*  
*SK Hynix*  
Icheon, South Korea  
chunseok.jeong@sk.com

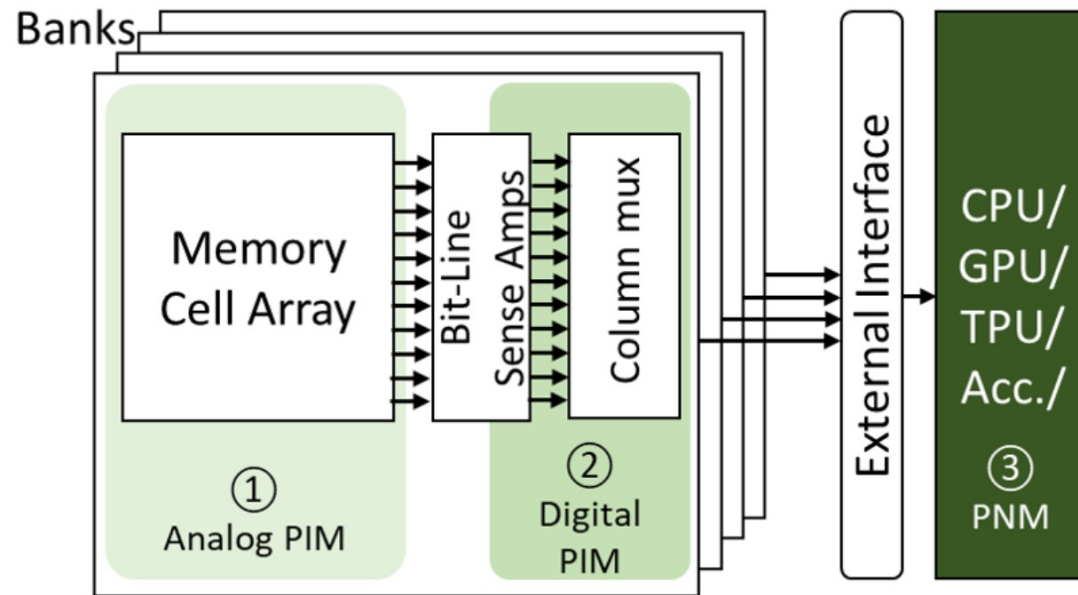
Seho Kim  
*DRAM Design*  
*SK Hynix*  
Icheon, South Korea  
seho5.kim@sk.com

Il Park  
*DRAM Design*  
*SK Hynix*  
Icheon, South Korea  
il.park@sk.com

Mithuna Thottethodi  
*Electrical and Computer Engineering*  
*Purdue University*  
West Lafayette, IN, U.S.A.  
mithuna@purdue.edu

T. N. Vijaykumar  
*Electrical and Computer Engineering*  
*Purdue University*  
West Lafayette, IN, U.S.A.  
vijay@ecn.purdue.edu

# Processing in/near Memory



## Analog PIM

- Low power, Fast
- **Accuracy loss**

## Digital PIM

- Low power, Fast
- **Limited PE**

## PNM

- Many computation units
- **Low bandwidth**

# When is digital PIM beneficial?

- **Compute-bound operation**
  - Requires large # of PEs and high data reuse
  - Not suitable for PIM
- **Memory-bound operation**

	small (high temporal reuse)	large (low temporal reuse)
small (high temporal reuse)	e.g. Tiled Conv	Matrix-vector multiplication
large (low temporal reuse)	Matrix-vector multiplication	Element-wise Operation (e.g. Add)

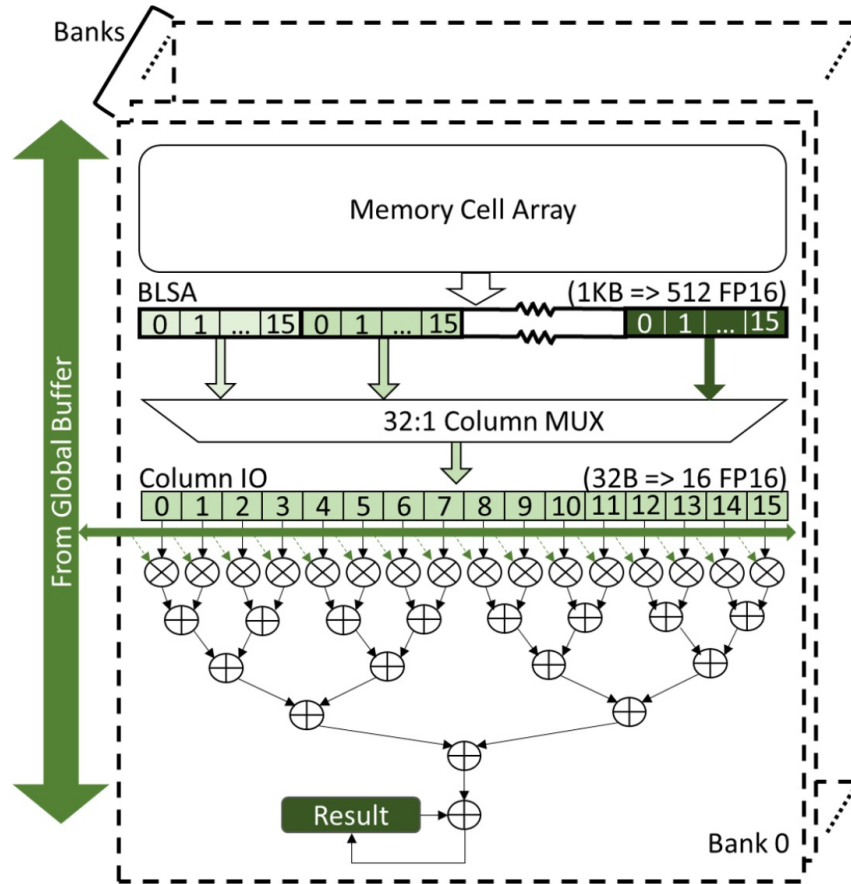
Matrix-Vector Multiplication

- LSTM, RNN, FC
- MAC operations

# Contributions

- Place a minimal compute of **only MAC units and buffers**
  - Previous works use an large number of PEs
    - e.g. superscalar, vector/SIMD
    - But PIM is subject to severe area and power constraints
- **DRAM command-like interface** for the host to issue commands
- Prevent PIM-host interface from becoming bottleneck
  - **gang** multiple compute operations
  - **complex** compute commands
  - targeted reduction of timing overhead (e.g.  $t_{FAW}$ )
- Reduce output vector write traffic
  - unusually-wide **interleaved layout for the filter matrix**

# Newton Datapath



All banks perform compute simultaneously (*gang*)

Reduction tree for reducing # MAC units

Interleaved accumulation for reducing output traffic

Assumption: BFloat16 / FP16

Justification: High accuracy for NLP task

# Newton's Tiled MV computation

---

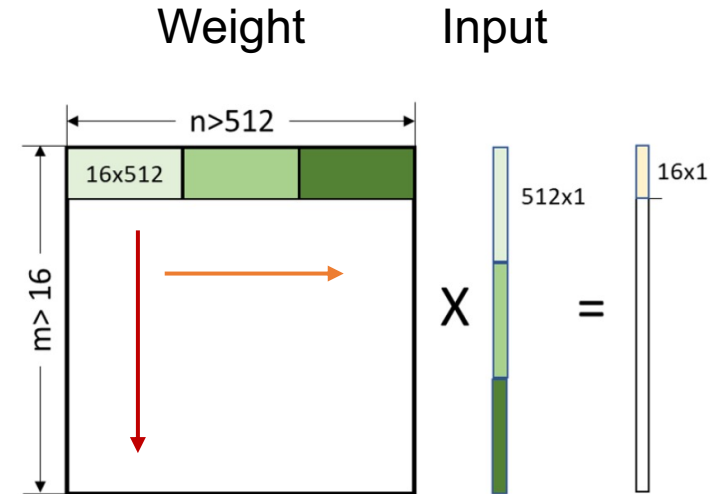
## Algorithm 1 Newton's Tiled MV computation

---

```

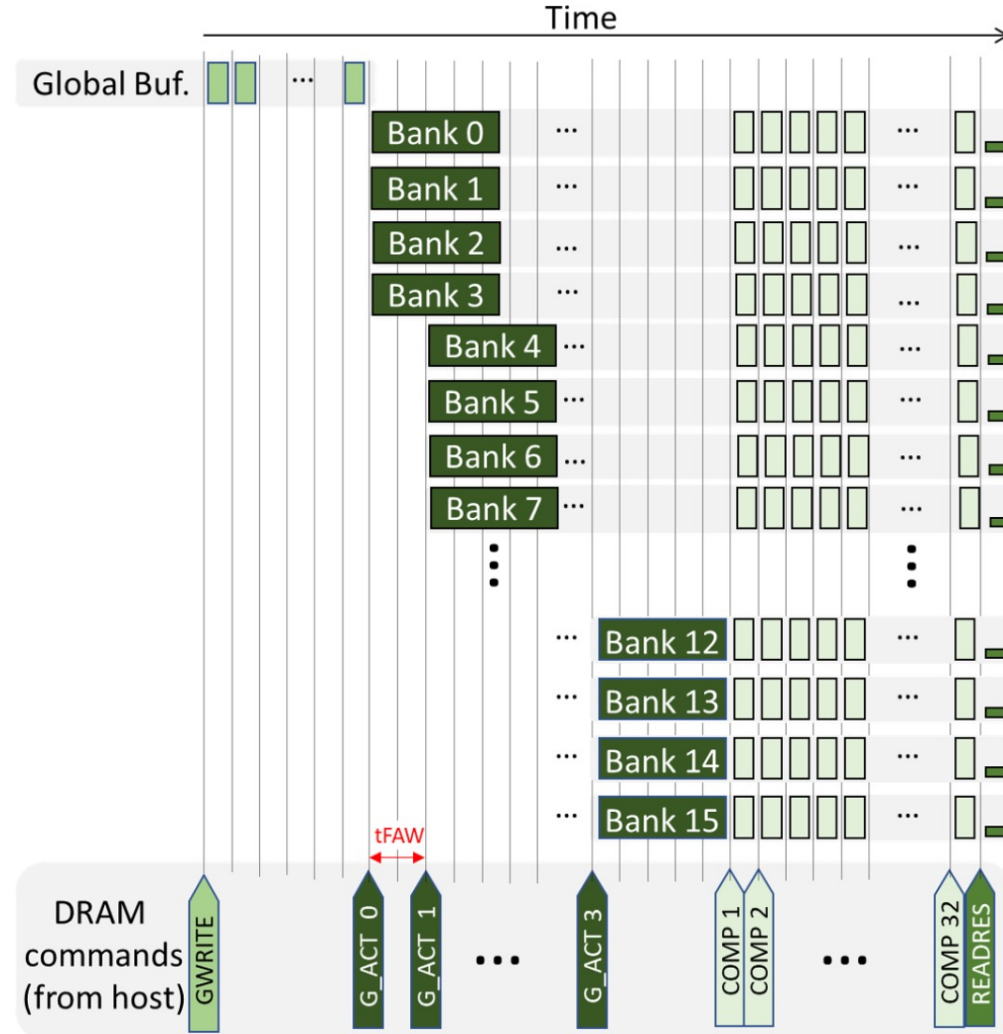
1: function MVPRODUCT(InputVector  $V$ , Matrix  $M$ ,  $m$ ,  $n$ )
2:    $numChunks = n/512$                                 ▷ Number of chunks
3:    $C[1..numChunks] \leftarrow split(V)$  ▷ Split vector to chunks
4:   for  $i \in 1..numChunks$  do                                ▷ Outermost loop.
5:      $GlobalBuffer \leftarrow C[i]$ 
6:      $r = m/16$                                 ▷ Number of vertical tile positions
7:     for  $j \in 1..r$  do
8:       ▷ Compute 1 DRAM row x Global Buffer per bank
9:       for all  $b \in 1..numBanks$  do
10:         $Results[b] \leftarrow ComputeTile(Tile\ j, Row\ b)$ 
11:      end for
12:       $TileResult \leftarrow ReadResultsFromAllBanks()$ 
13:      ▷ Tile result sent for accumulation at host
14:    end for
15:  end function
  
```

---



Interleaved (row-first) to reduce output buffer traffic

# Newton computation



Command	Operation
COMP#	Ganged multiply of sub-chunk# in all banks
READRES	Read the Result latches of all banks
GWRITE#	WRITE sub-chunk# to the Global Buffer
G_ACT#	Ganged activation of 4-bank cluster#

*Complex command → low PIM-host BW consumption*

# Methodology

- Benchmarks

Workload	Matrix	Vector
GNMT LSTMs1 [45]	$4096 \times 1024$	$1024 \times 1$
GNMT LSTMs2 [45]	$4096 \times 2048$	$2048 \times 1$
BERTs1 [10]	$1024 \times 1024$	$1024 \times 1$
BERTs2 [10]	$1024 \times 4096$	$4096 \times 1$
BERTs3 [10]	$4096 \times 1024$	$1024 \times 1$
AlexNetL6 [26]	$21632 \times 2048$	$2048 \times 1$
AlexNetL7 [26]	$2048 \times 2048$	$2048 \times 1$
DLRMs1 [31]	$512 \times 256$	$256 \times 1$

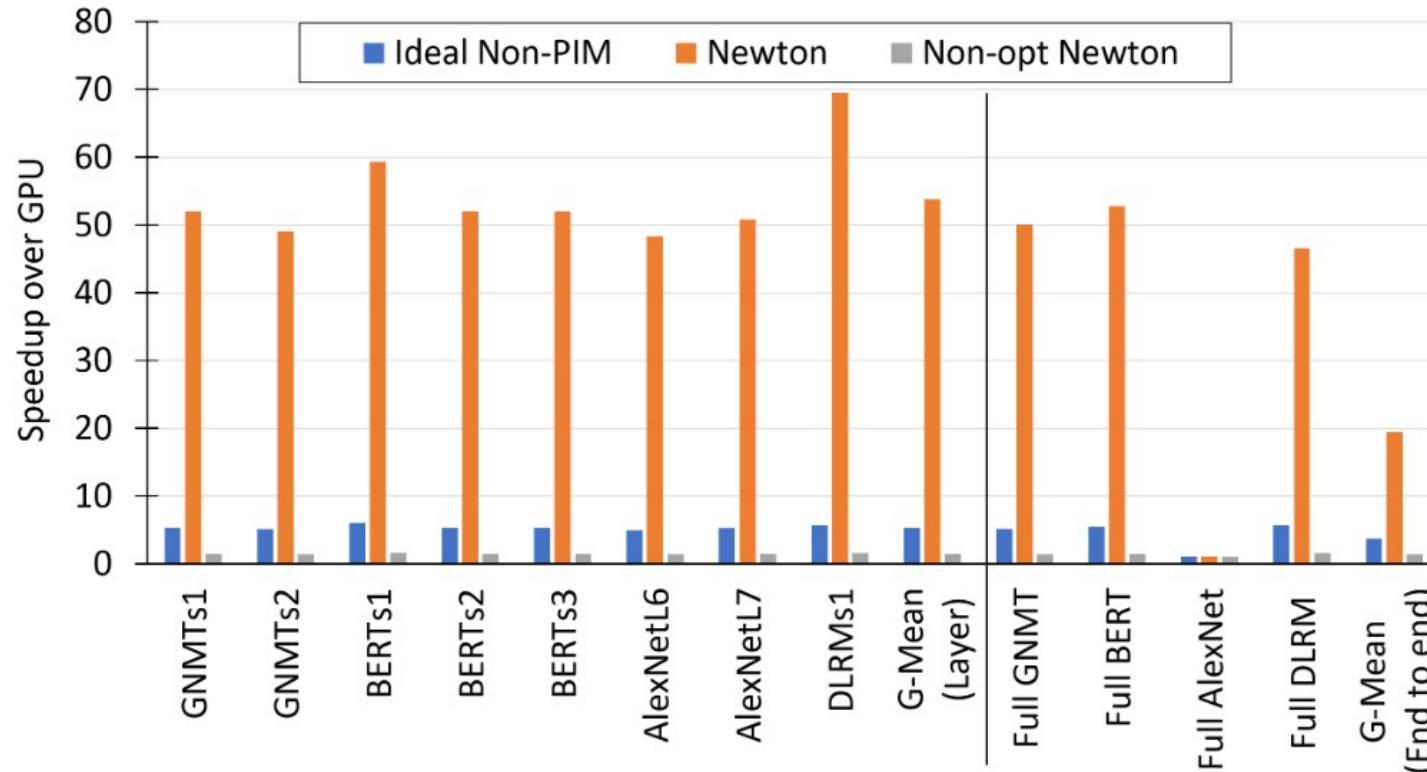
- DRAM configuration

Num of Ranks	1
Num of Banks	16
Num of Rows in each bank	32768
Num of Column I/Os per row	32
Column I/O bit width	256b (16 bfloat16)
Num of Multipliers per bank	16
<b>Timing Parameters (in nanoseconds)</b>	
$t_{AA} = 22\text{--}29$ ns; $t_{RP} = 14$ ns; $t_{RCD} = 14$ ns; $t_{RAS} = 33$ ns	

Only show partial parameters (proprietary)  
No power parameters (proprietary)

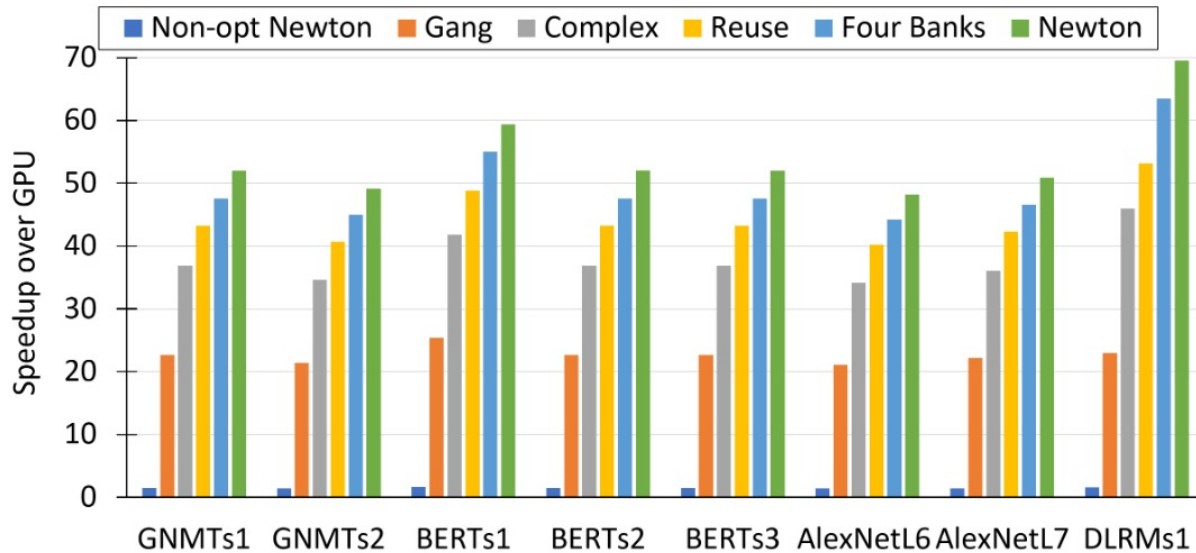


# Evaluation

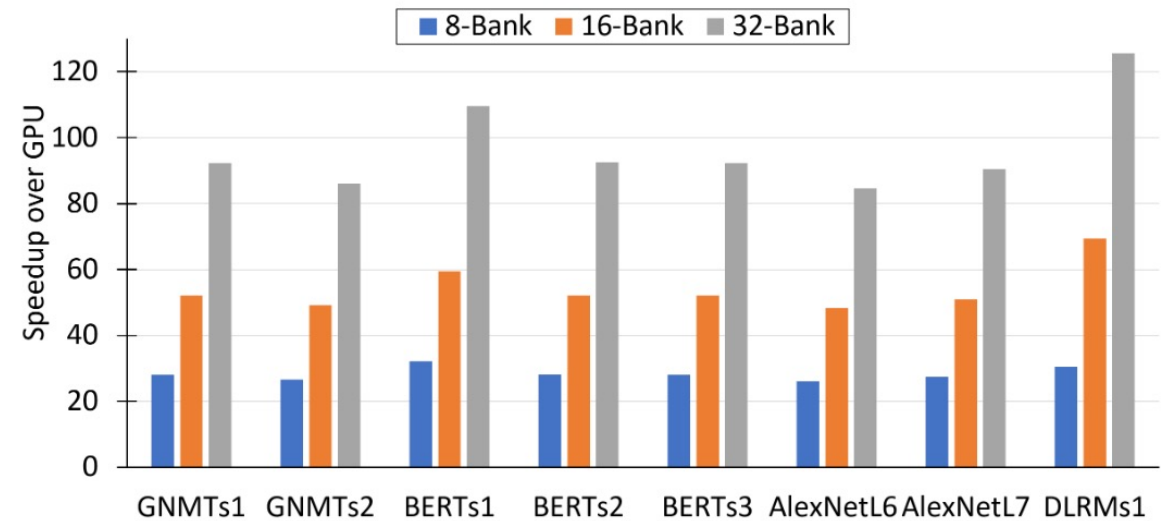


- Newton's speedup comes from **BW reduction**, not from computation capability (vs Ideal Non-PIM)
- DLRMs1 shows better performance since there's no DRAM refresh (short computation)

# Sensitivity Study

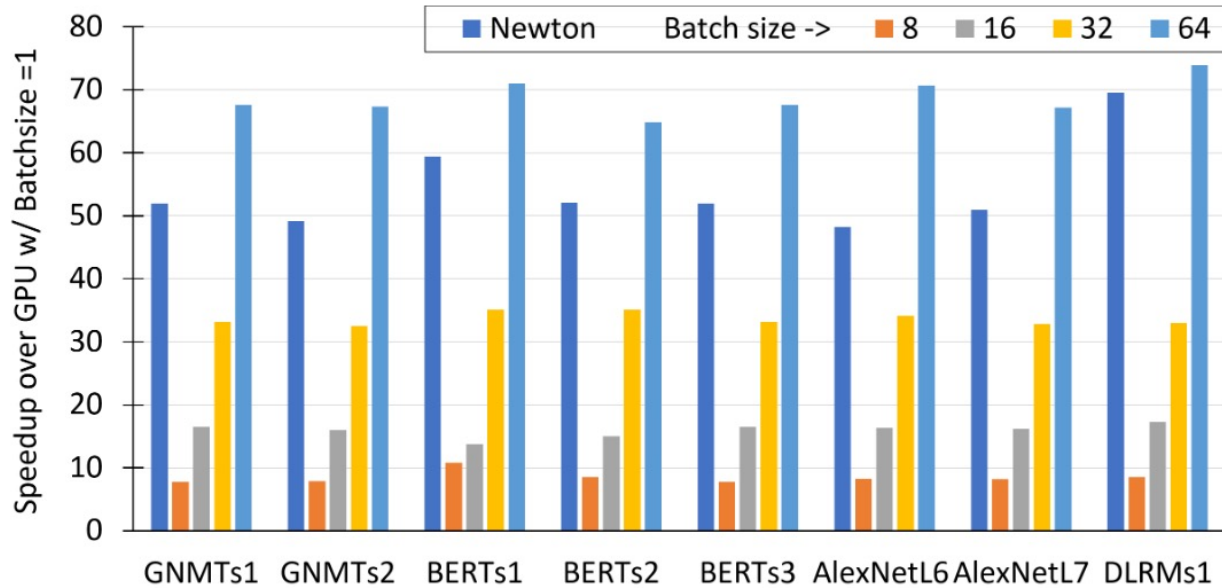


- Gang: all-bank ganged compute commands
- Complex: multi-step compute commands
- Reuse: reuse via tiling and interleaved layout for the filter matrix
- Four Banks: four-bank ganged activations
- Newton: aggressive  $t_{FAW}$



- Nearly perfect scaling

# Sensitivity Study (Batch size)



- Newton is good for inference (batch size  $\leq 8$ )
- More data-reuse in GPU  $\rightarrow$  Less profitability for Newton

# Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology

ISCA'21

Industrial Product

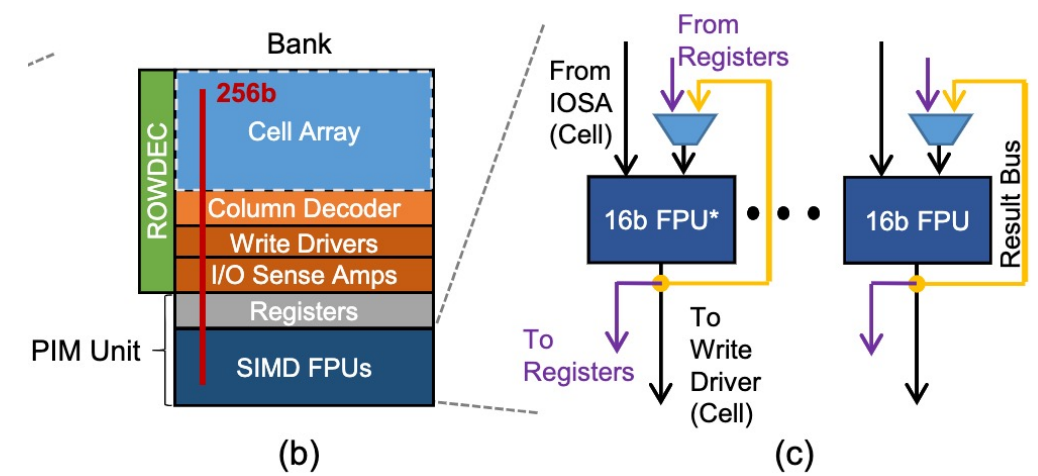
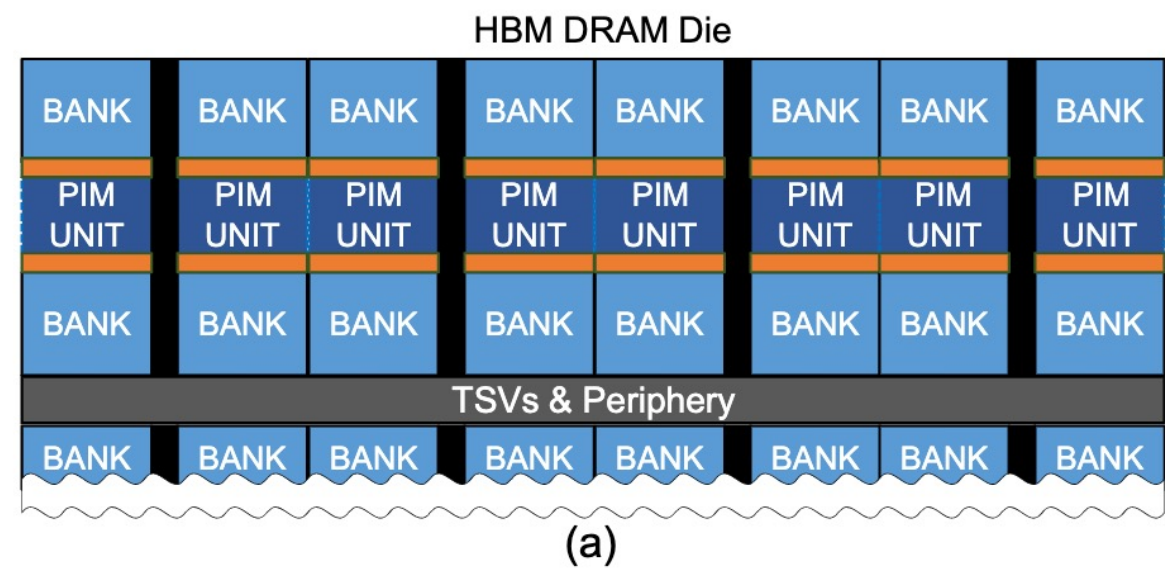
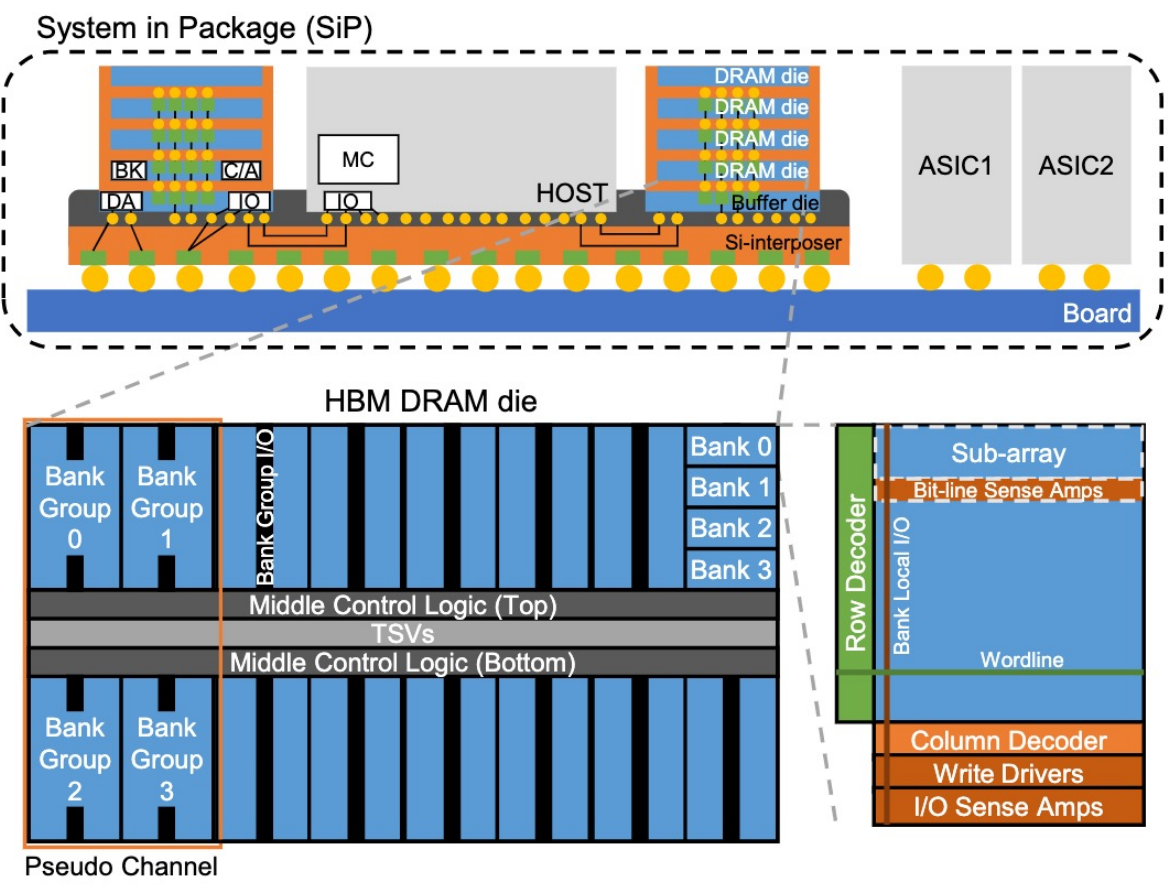
Sukhan Lee<sup>§1</sup>, Shin-haeng Kang<sup>§1</sup>, Jaehoon Lee<sup>1</sup>, Hyeonsu Kim<sup>2</sup>, Eojin Lee<sup>1</sup>, Seungwoo Seo<sup>2</sup>,  
Hosang Yoon<sup>2</sup>, Seungwon Lee<sup>2</sup>, Kyoungwan Lim<sup>1</sup>, Hyunsung Shin<sup>1</sup>, Jinhyun Kim<sup>1</sup>,  
Seongil O<sup>1</sup>, Anand Iyer<sup>3</sup>, David Wang<sup>3</sup>, Kyomin Sohn<sup>1</sup> and Nam Sung Kim<sup>§1</sup>

<sup>1</sup>Memory Business Division, Samsung Electronics

<sup>2</sup>Samsung Advanced Institute of Technology, Samsung Electronics

<sup>3</sup>Device Solutions America, Samsung Electronics

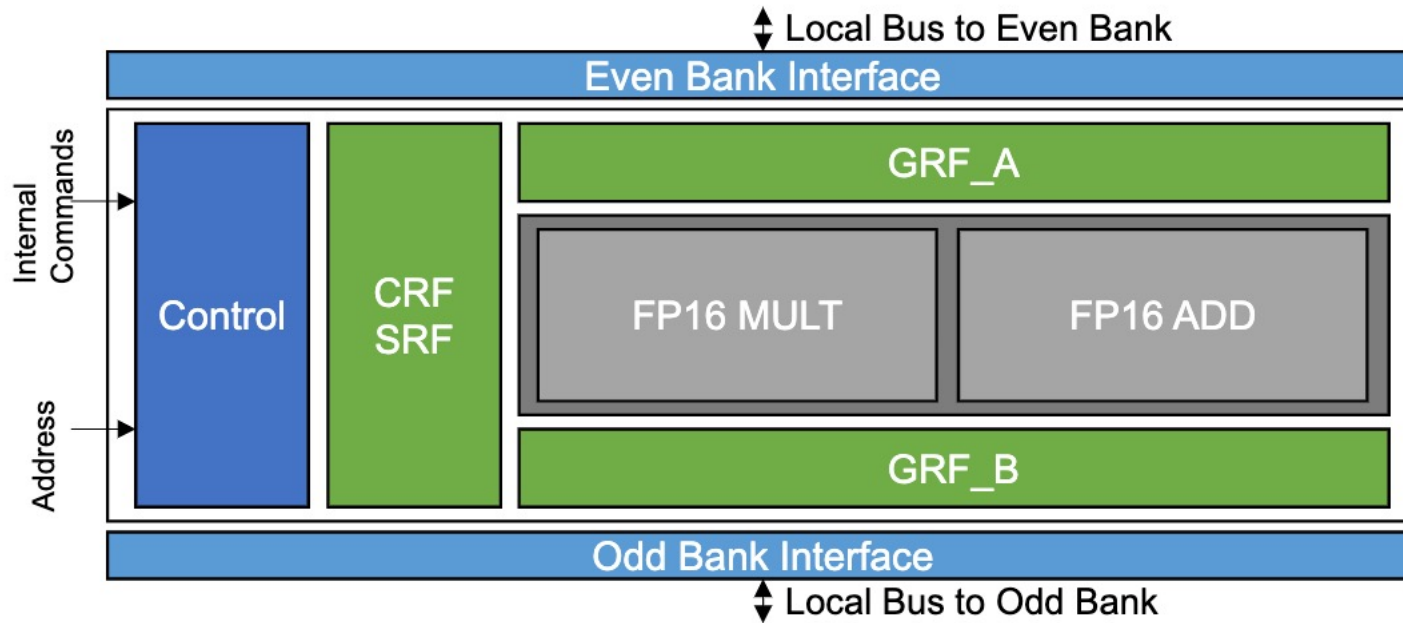
# HBM structure and PIM units



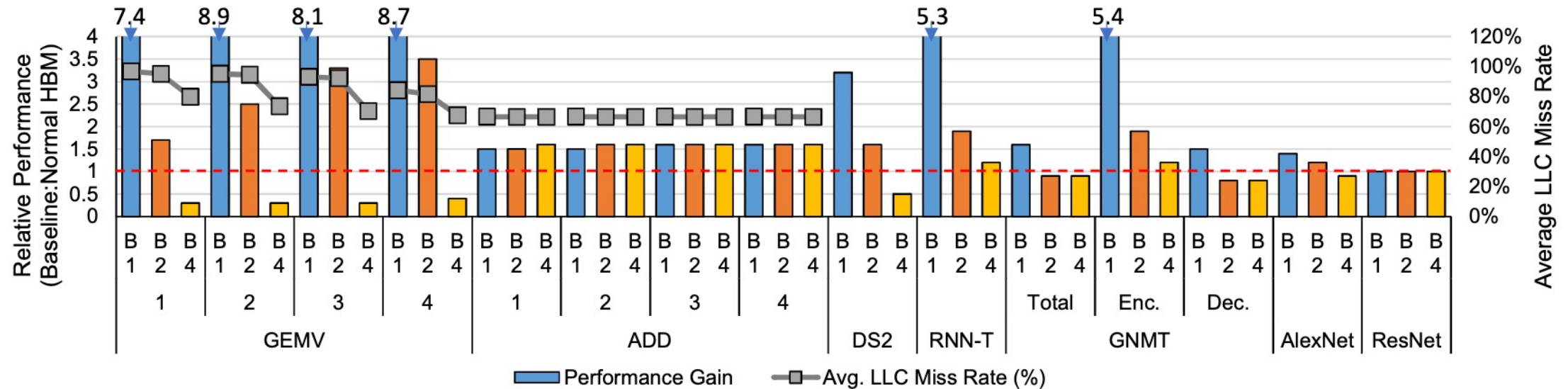
# RISC-V-like instruction format

TABLE III: The PIM-HBM instruction format.

	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
Control	OPCODE				U									IMM0							IMM1											
Data	OPCODE				DST			SRC0			U							R	U	DST #		U	SRC0 #		U	SRC1 #						
ALU	OPCODE				DST			SRC0			SRC1			SRC2		A	U		U	DST #		U	SRC0 #		U	SRC1 #						



# Evaluation



- Much less speedup than Newton (IF/ID and control overhead,  $\frac{1}{4}$  PEs)
- Support wide range of operations (e.g. Add, BN)

# Drawbacks of digital PIM

- Tricky cache coherence (Both)
- Lack of software stack (Newton)
  - Samsung is building its own stack
- Limited operation support (Newton)
  - What about other DL models?