

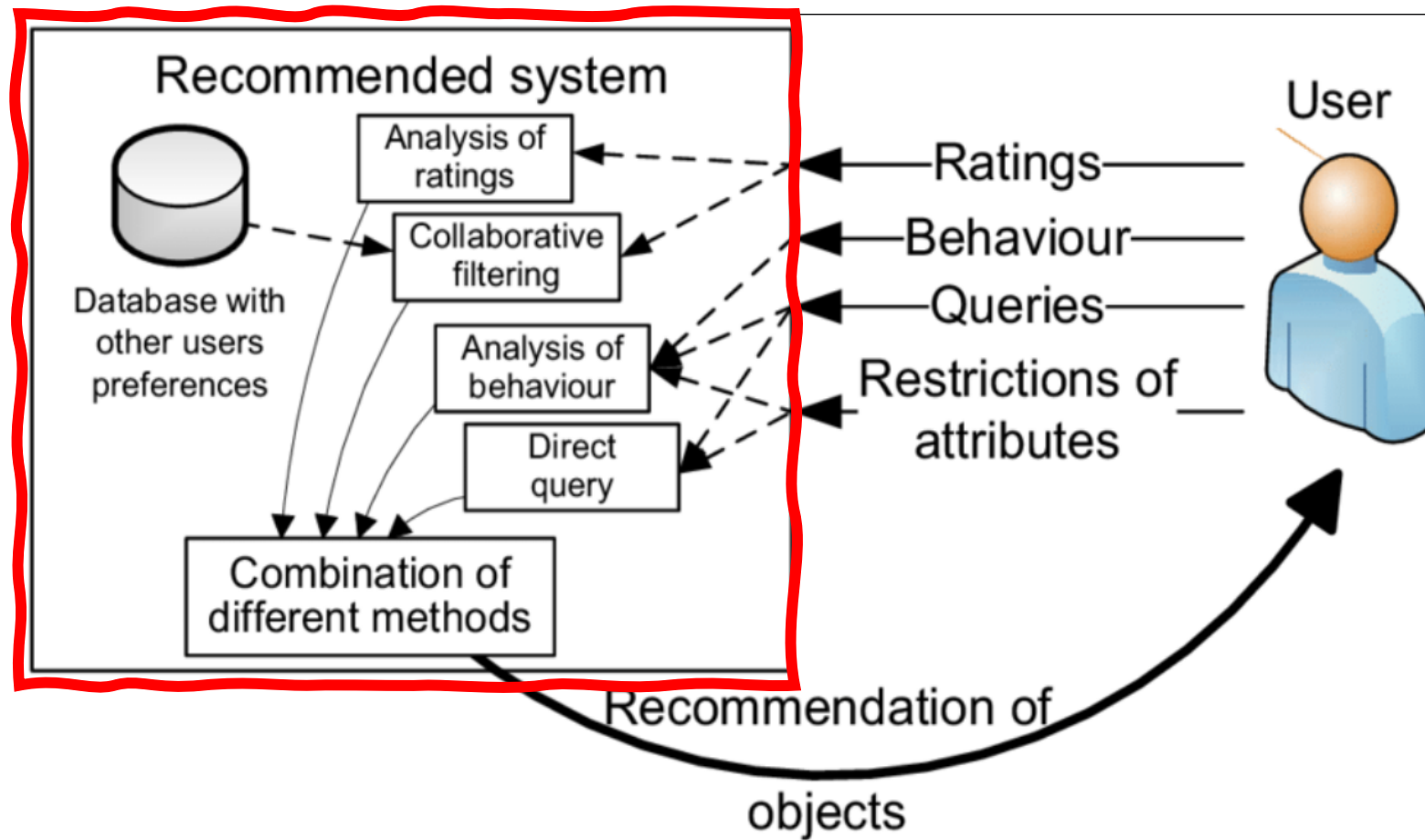
# PIM-SUM: PopCount-based Fast and Reliable In-memory Summation Scheme

*Fan Li, Ruizhi Zhu, Huize Li, Di Wu, Xin Xin*

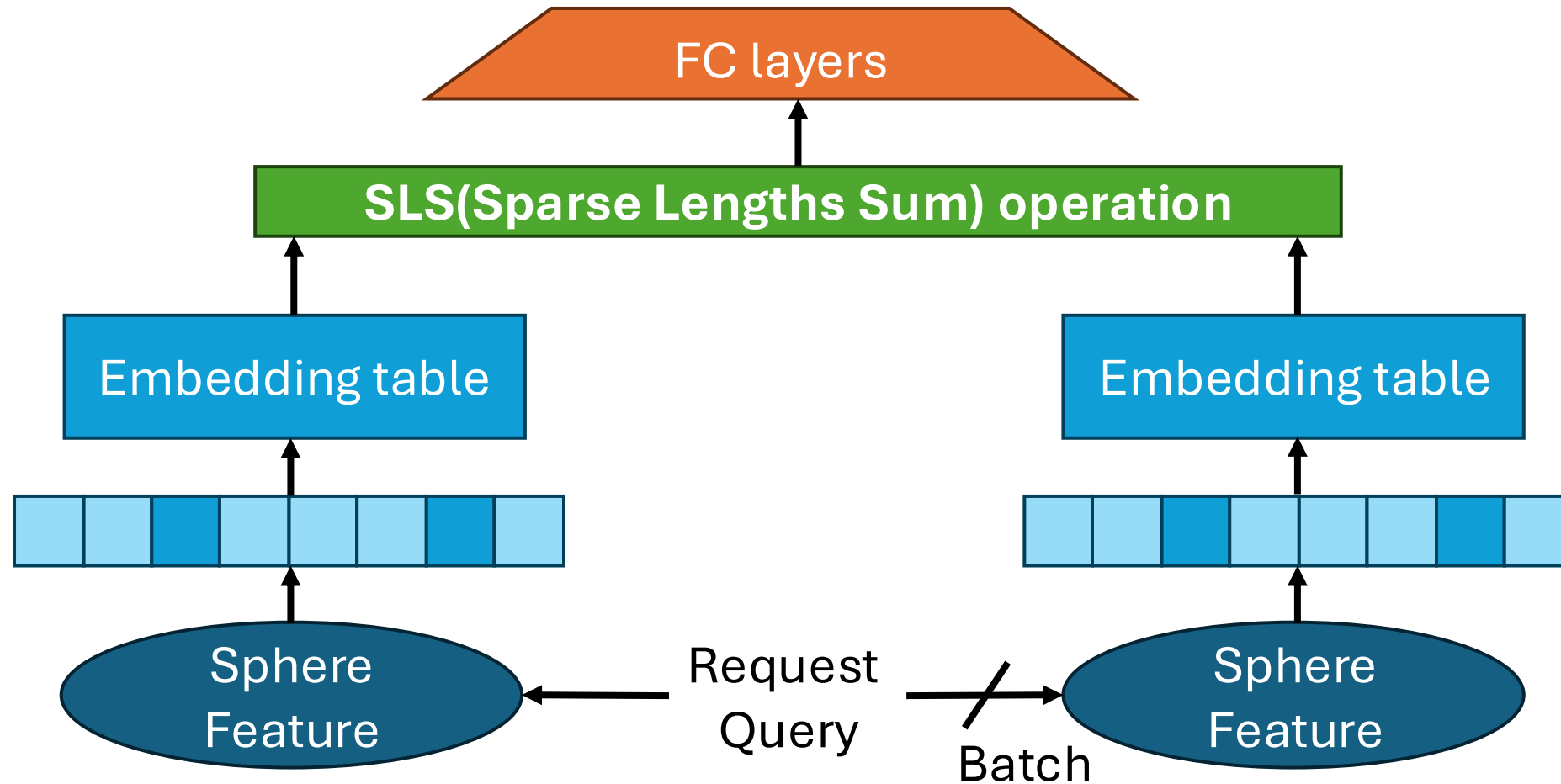
College of Engineering and Computer Science

University of Central Florida

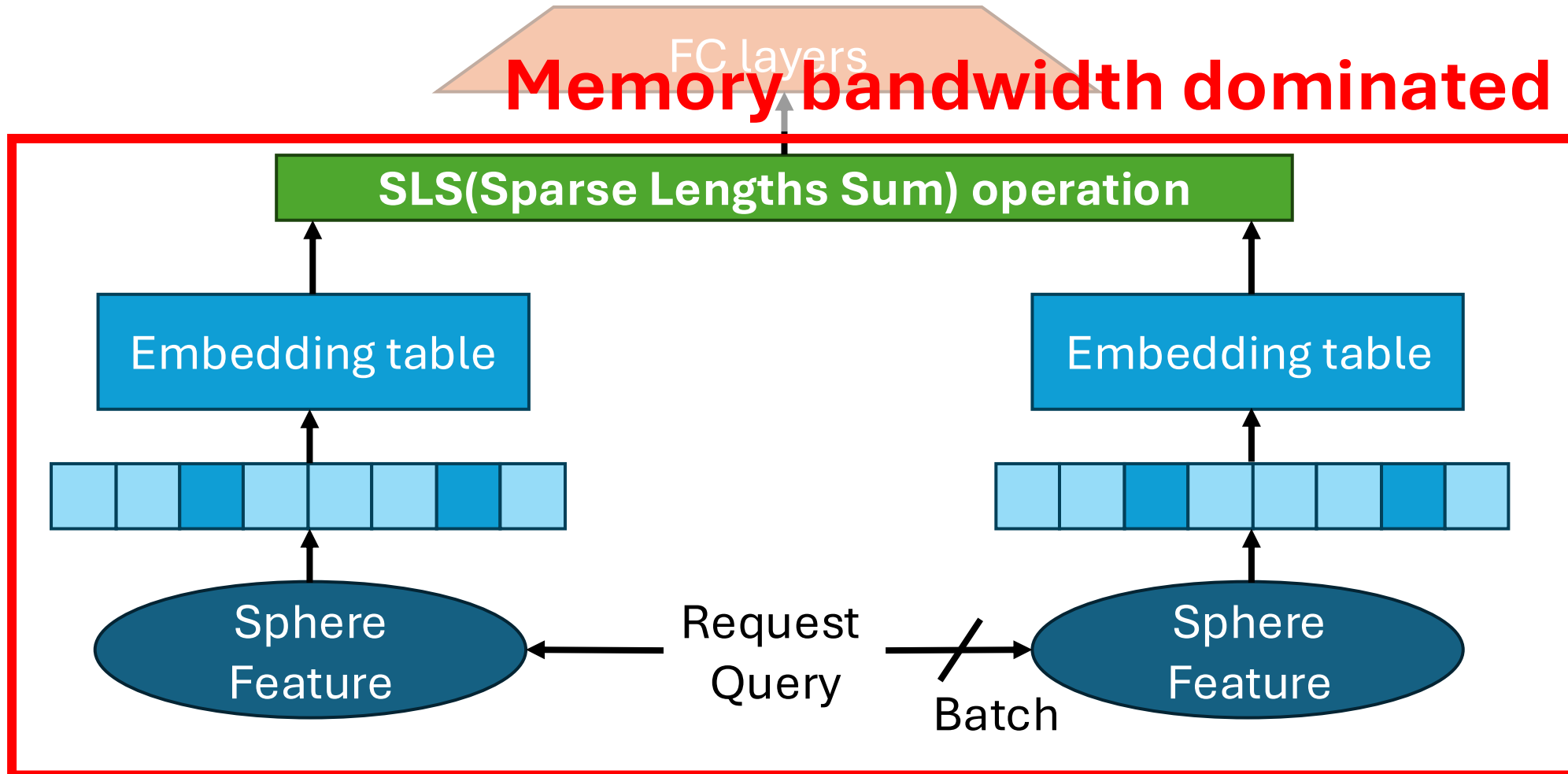
# Motivation--Recommendation System



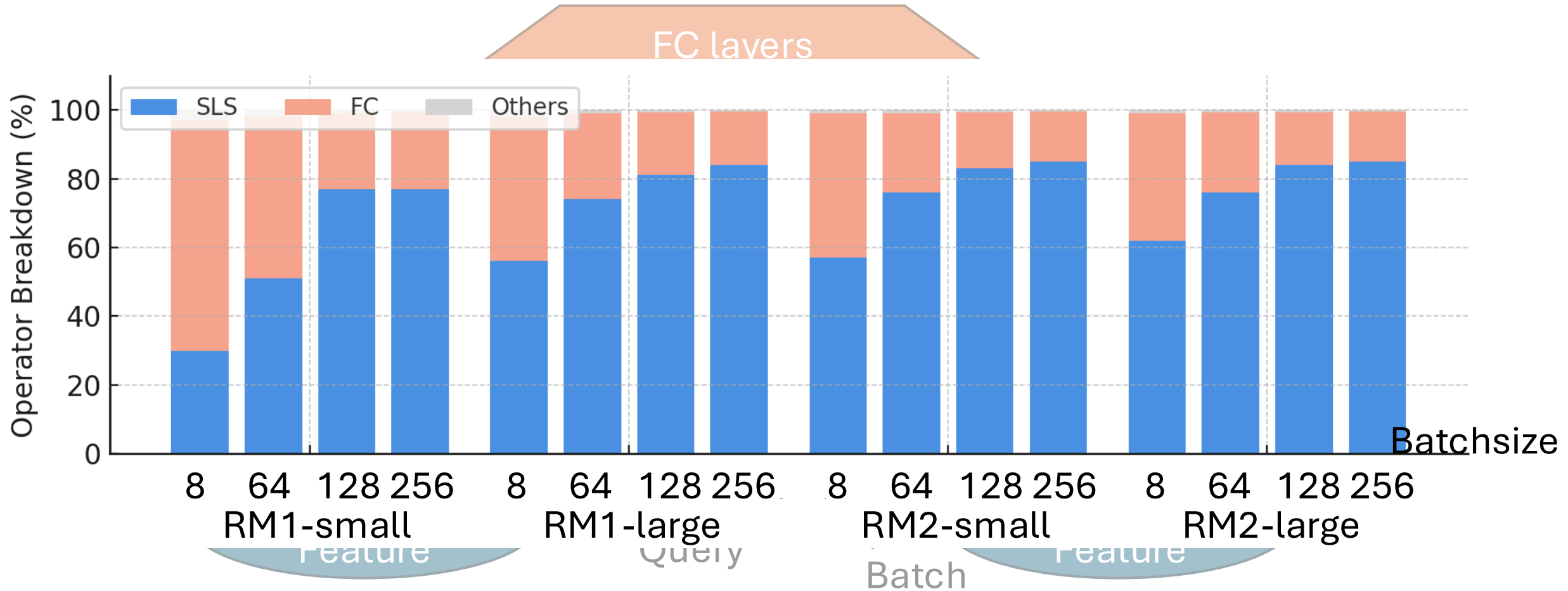
# Bottleneck of Recommendation Model



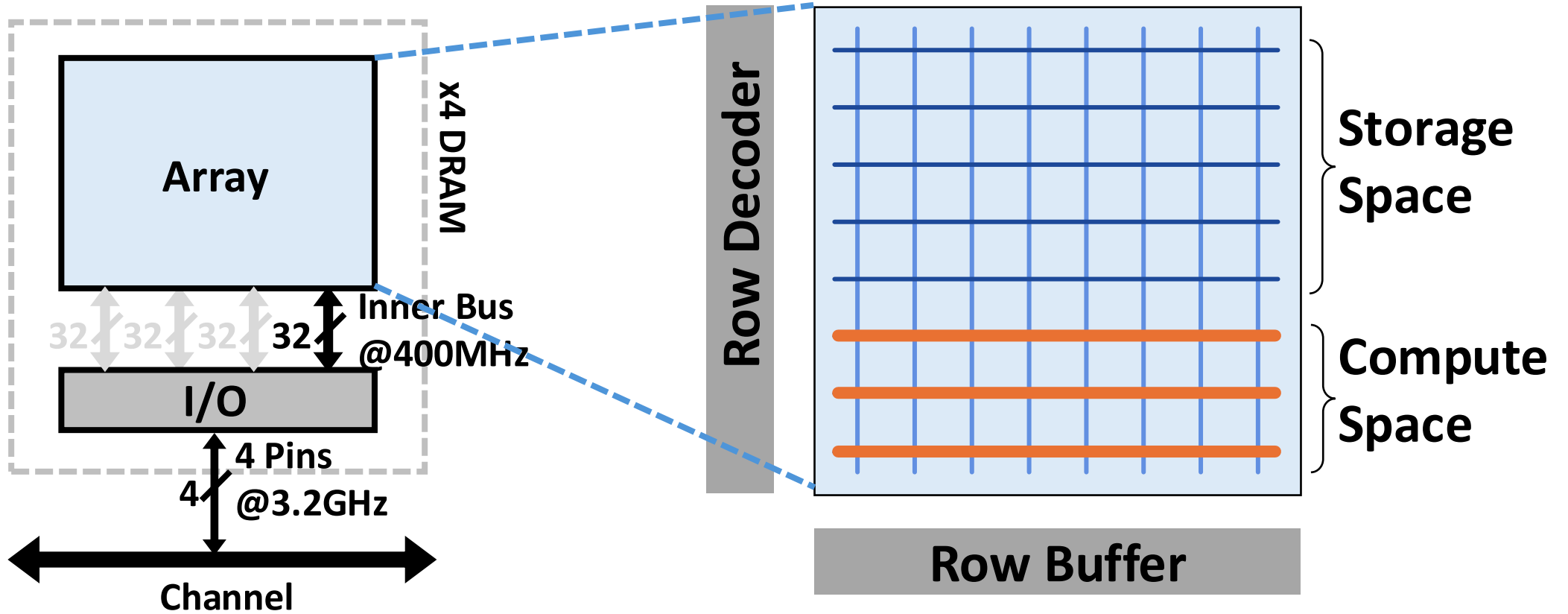
# Bottleneck of Recommendation Model



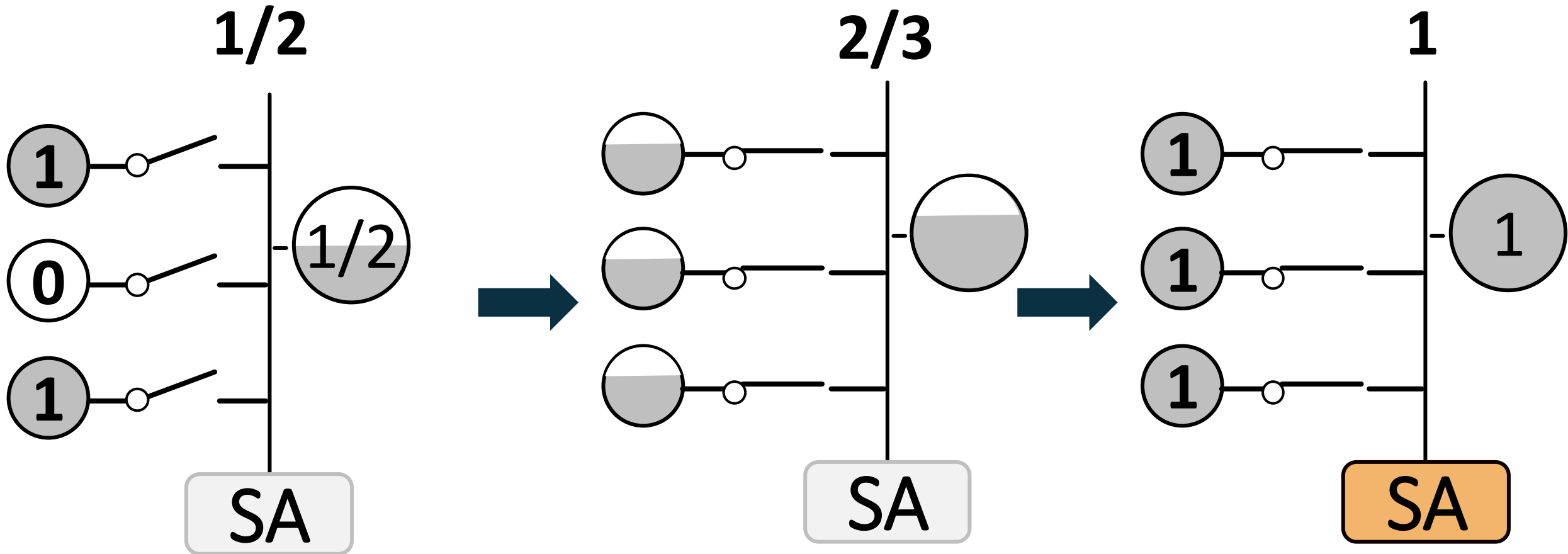
# Bottleneck of Recommendation Model



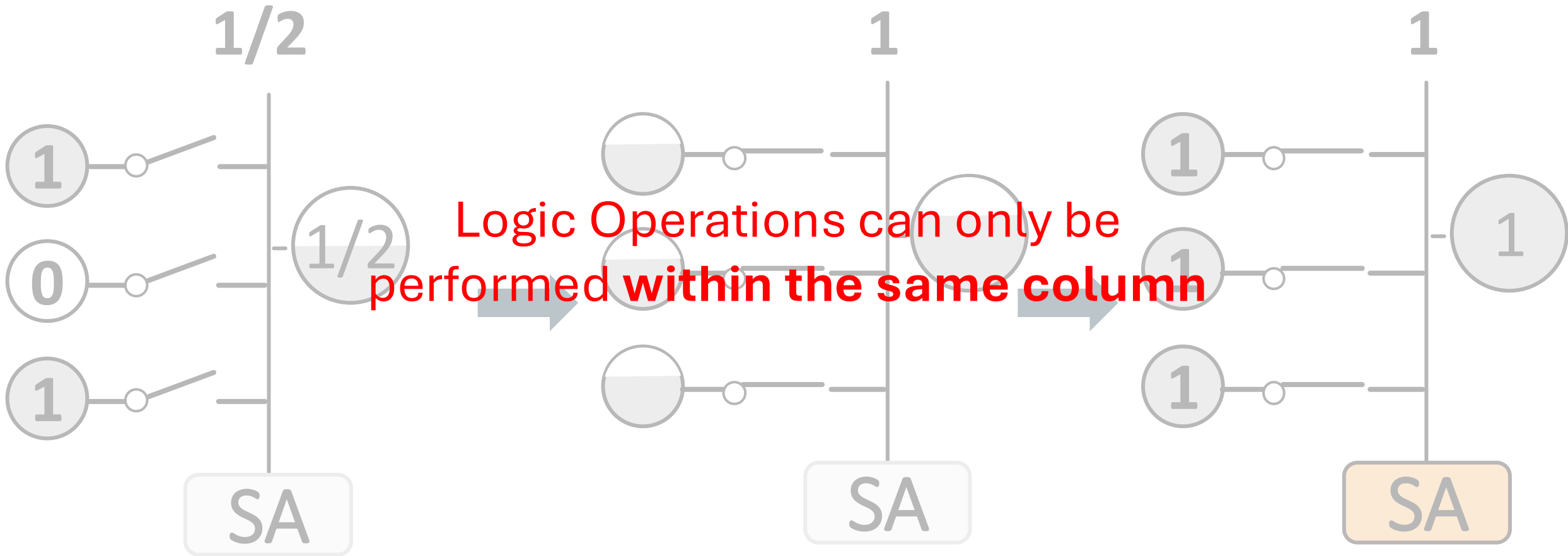
# Background—Processing in Memory(PIM)



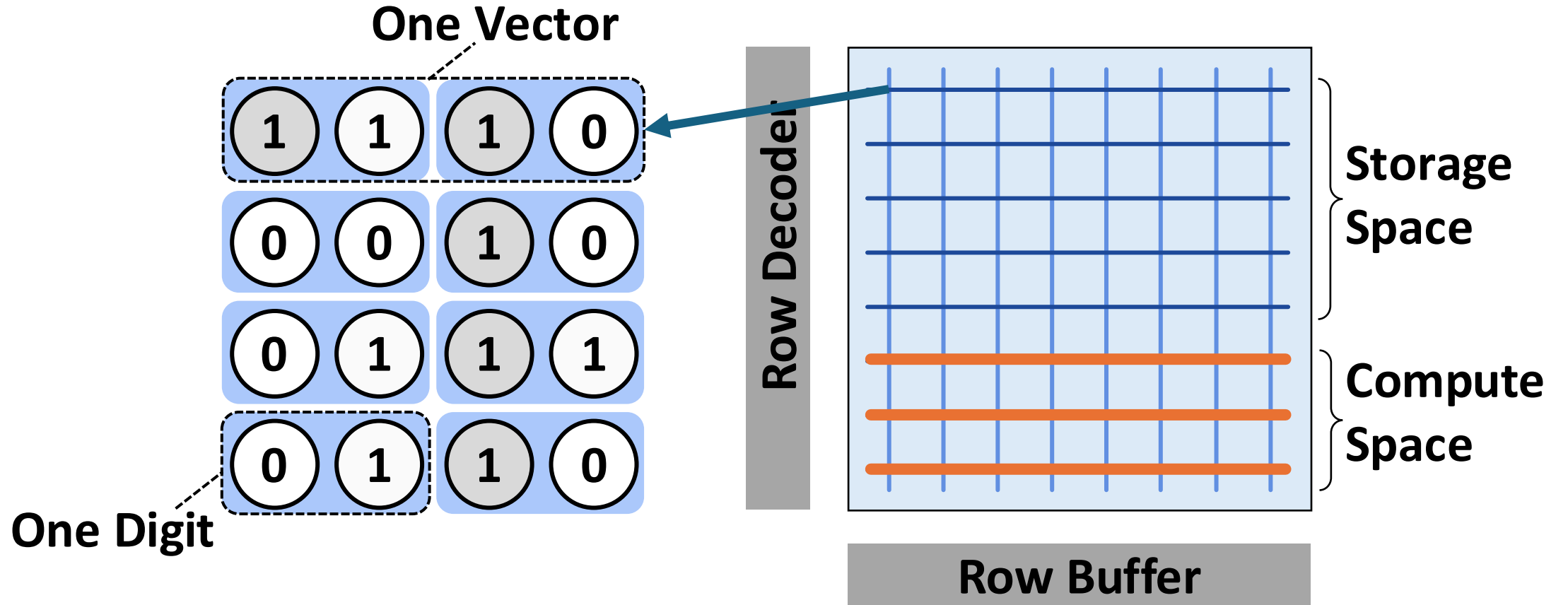
# Background—Processing in Memory(PIM)



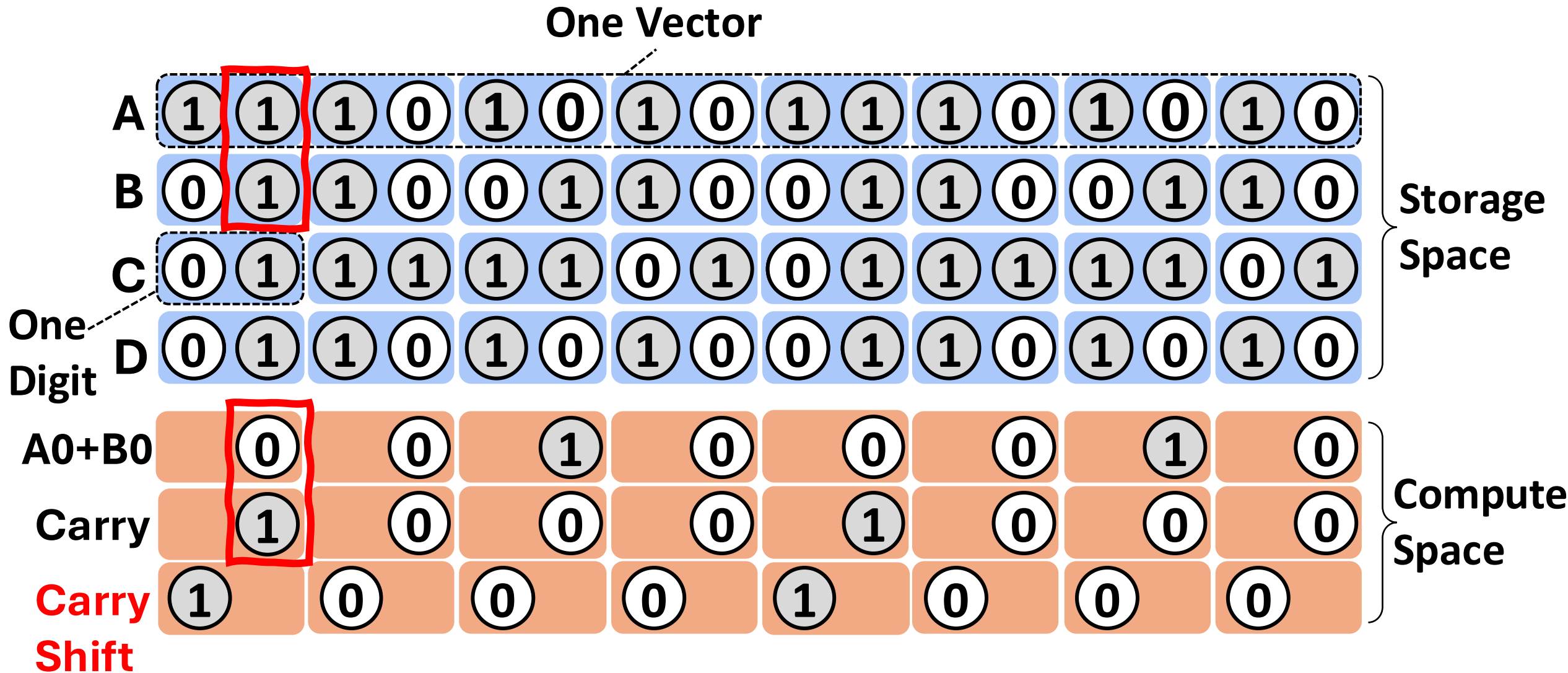
# Background—Processing in Memory(PIM)



# Mapping

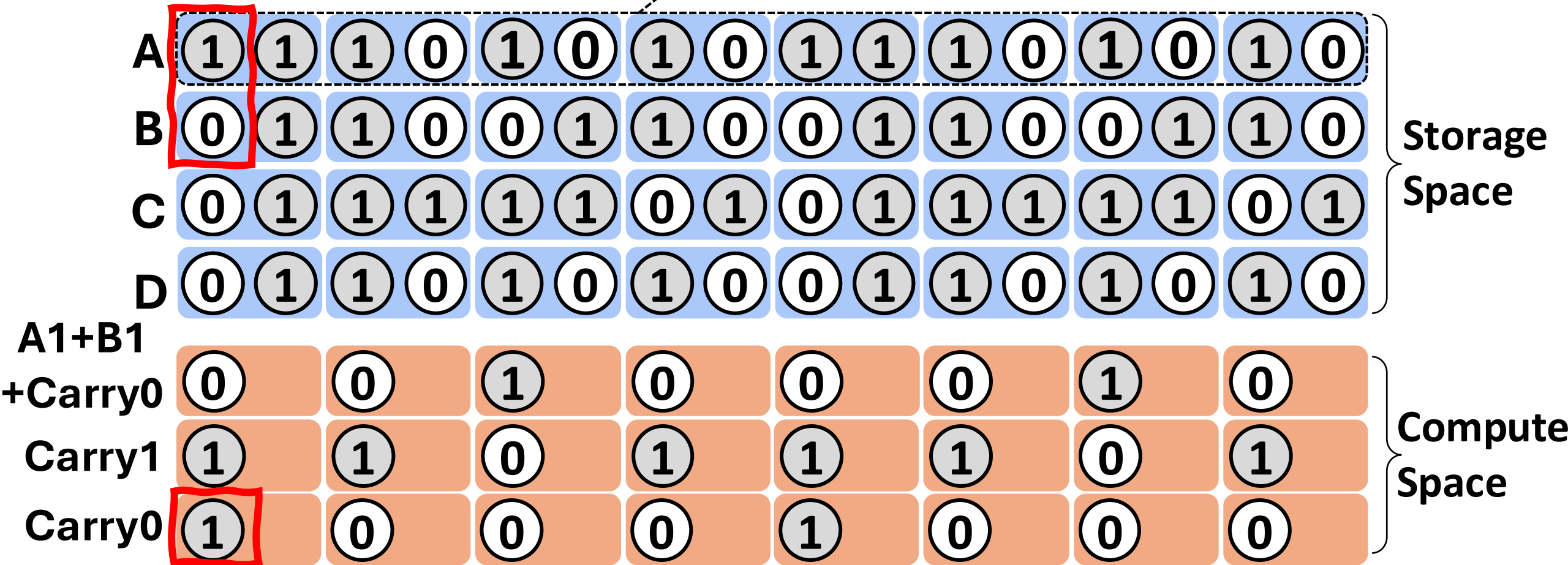


# Baseline



# Baseline

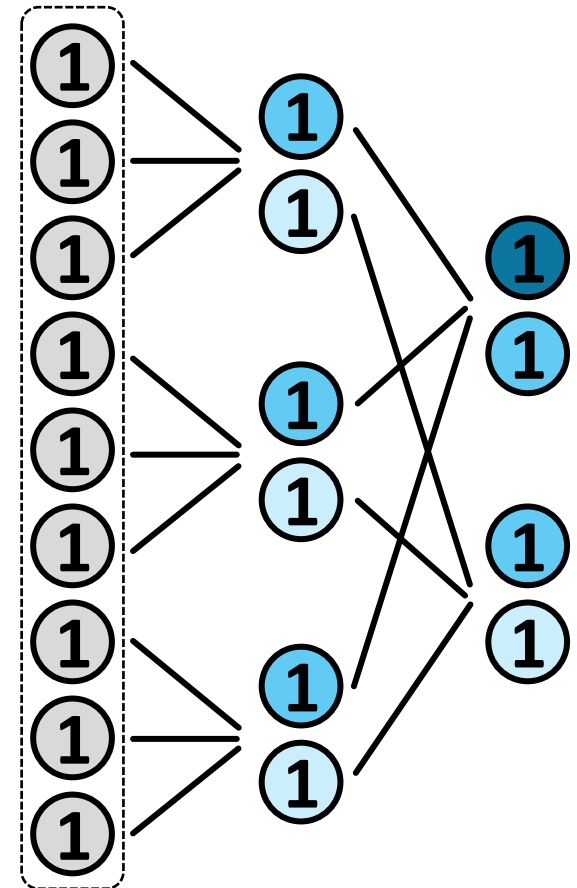
One Vector



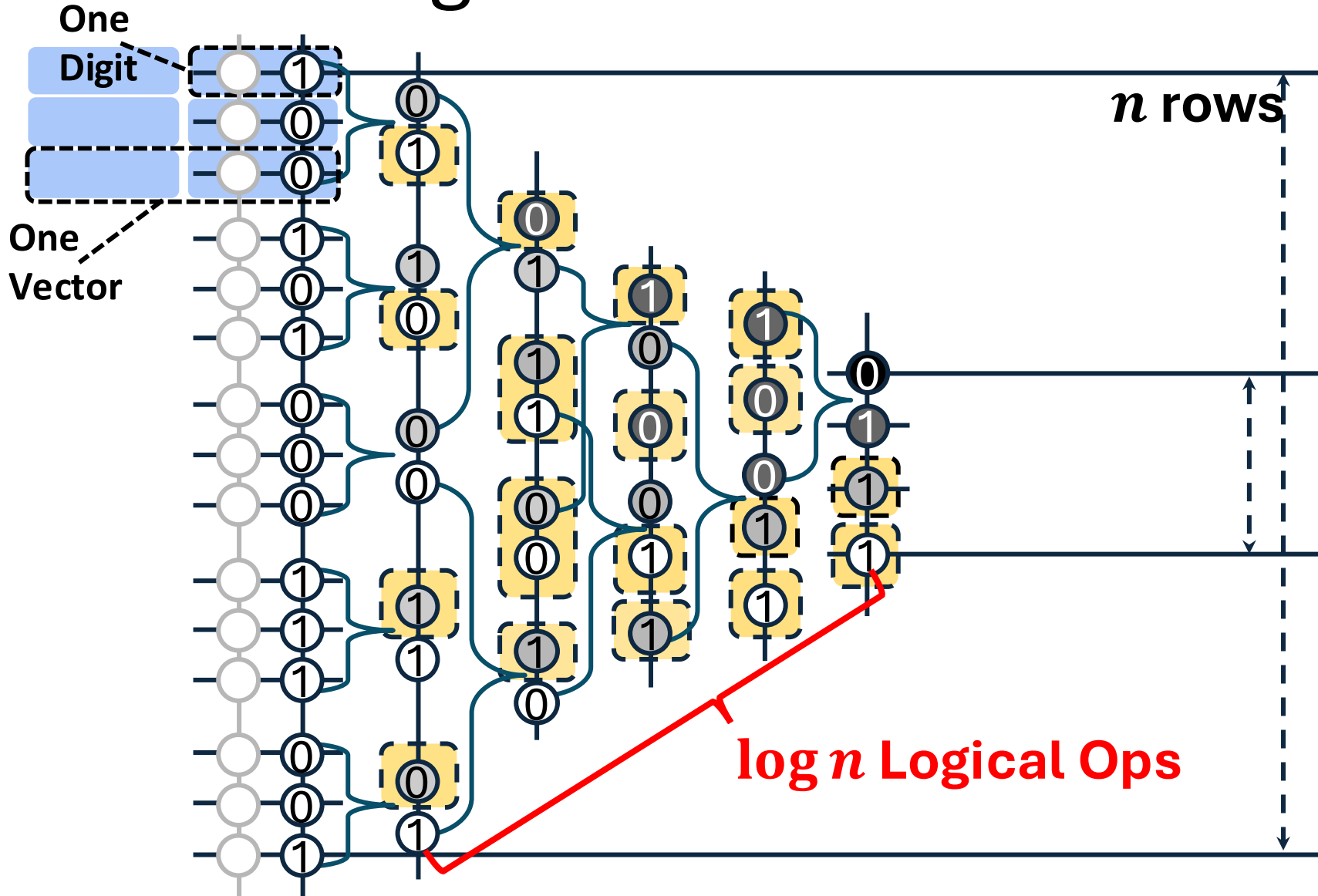
# Our Design--PIMSUM

## ***Steps:***

1. Aggregate 3 rows in the same level into 2 rows on the next level.
2. Pass unaggregated rows into future process



# Our Design--PIMSUM



Waiting for next stage

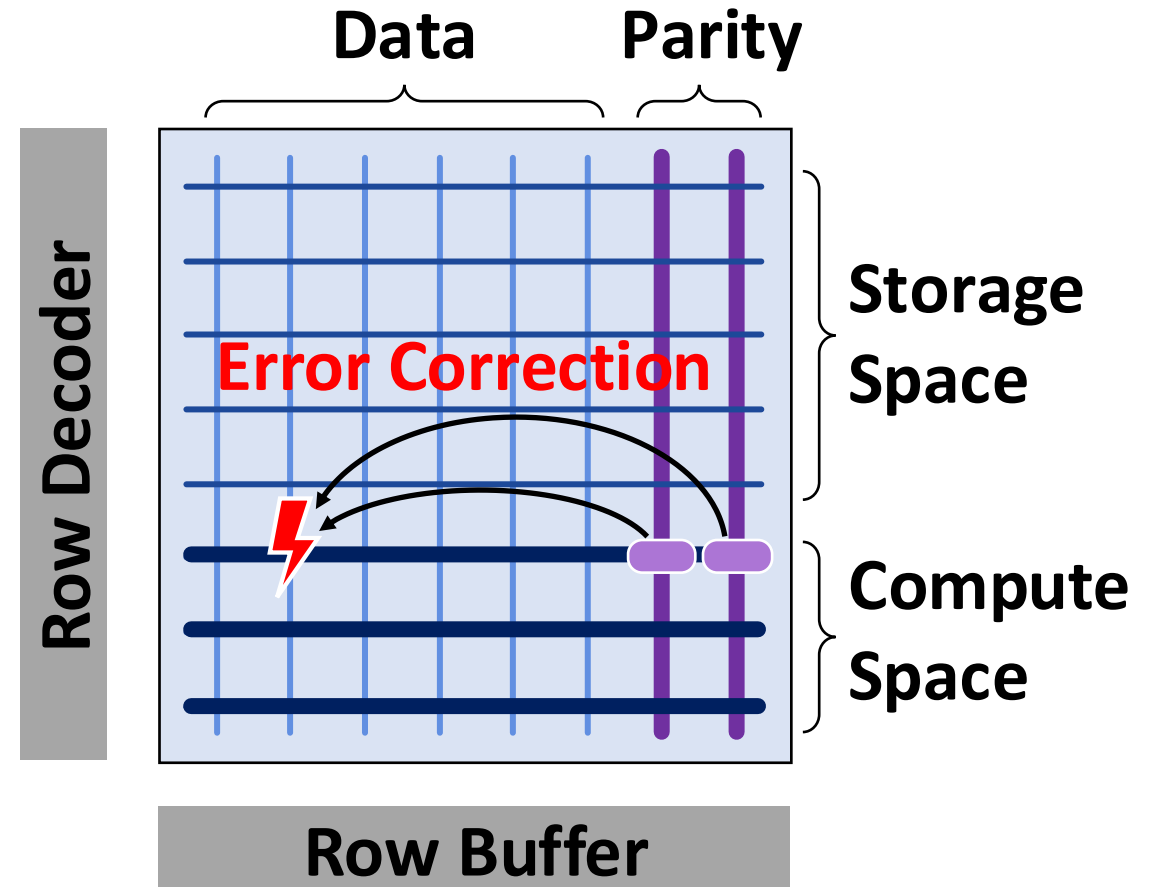


# Reliability--PIMSUM

$$\begin{cases} p_{a0} = a_0 + a_1 + a_2 \\ p_{a1} = a_0 + 2a_1 + 3a_2 \end{cases}$$



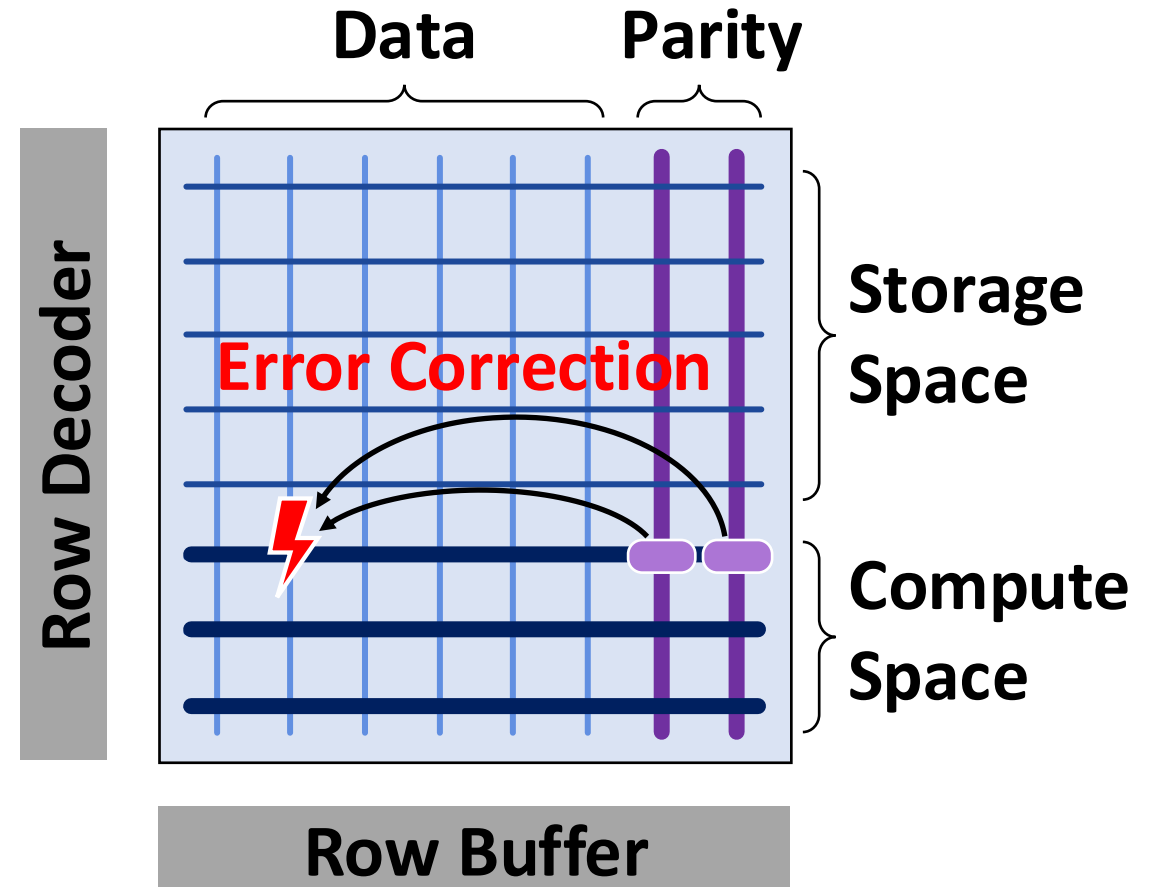
$$\begin{cases} p_{a0} = a_0 + (a_1 + e) + a_2 \\ p_{a1} = a_0 + 2(a_1 + e) + 3a_2 \end{cases}$$



# Reliability--PIMSUM

$$\begin{cases} p_{a0} = a_0 + a_1 + a_2 \\ p_{a1} = a_0 + 2a_1 + 3a_2 \end{cases}$$

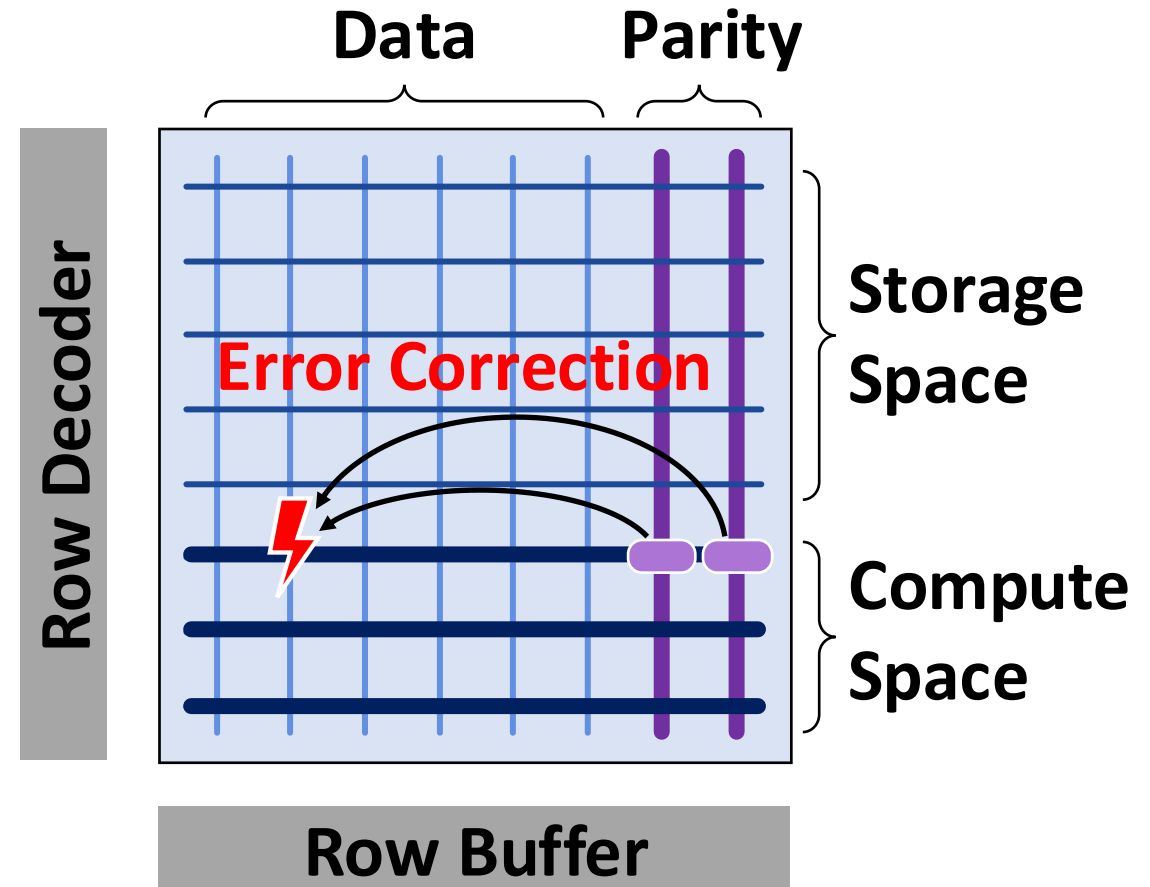
$$\begin{cases} p_{a0} = a_0 + (a_1 + e) + a_2 \\ p_{a1} = a_0 + 2a_1 + 2e + 3a_2 \end{cases}$$



# Reliability--PIMSUM

$$e_{ind} = \frac{p_{a1}' - p_{a1}}{p_{a0}' - p_{a0}}$$

$$e = p_{a0}' - p_{a0}$$

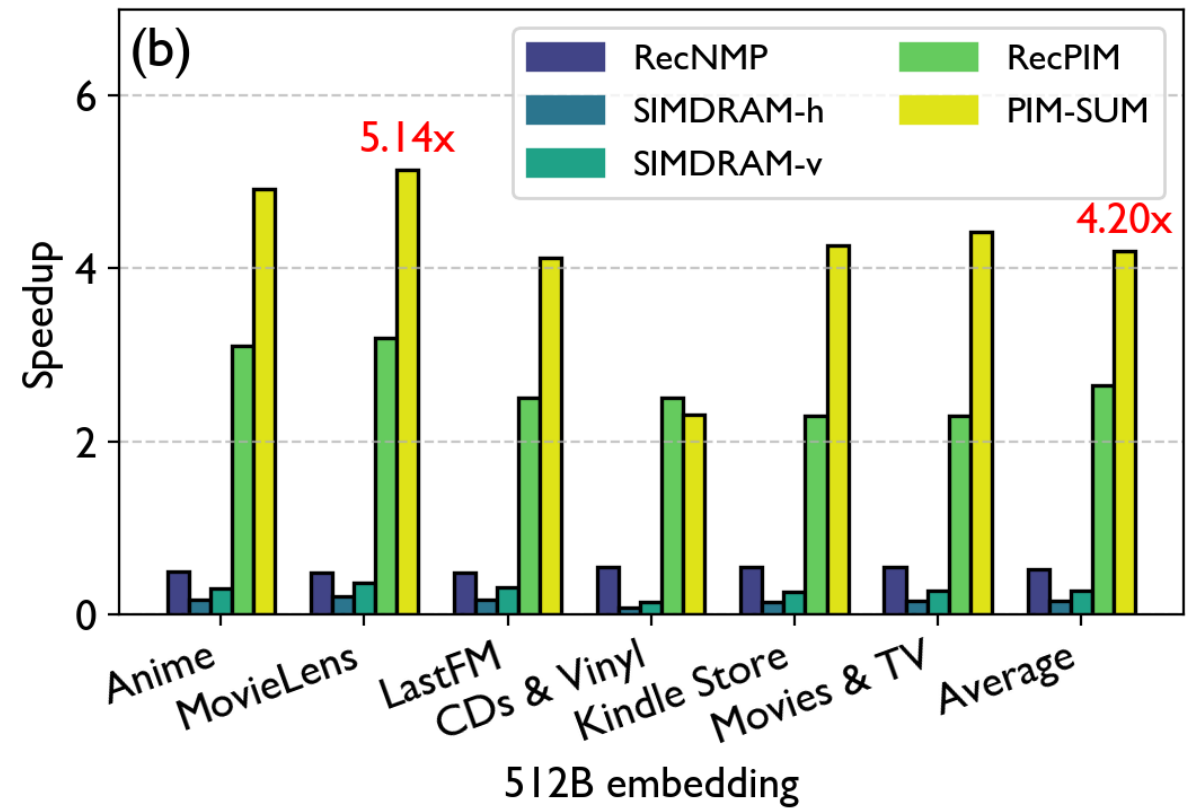
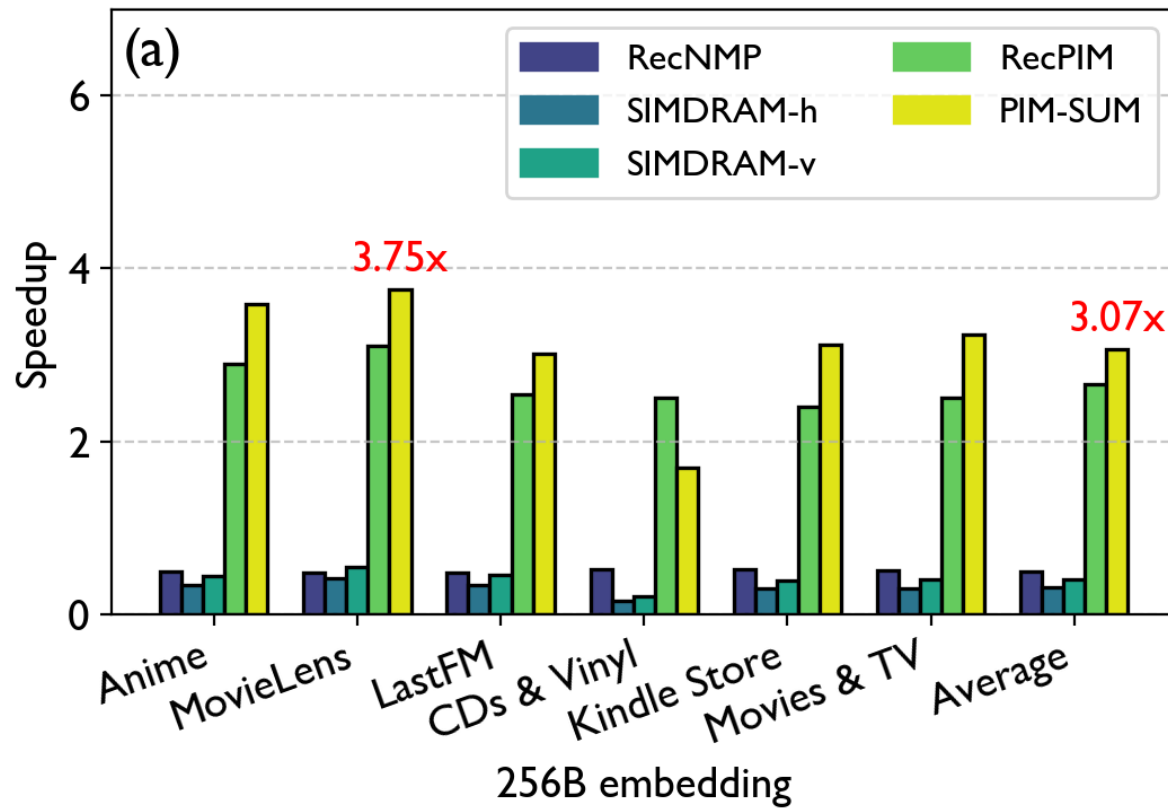


# Summary--PIMSUM

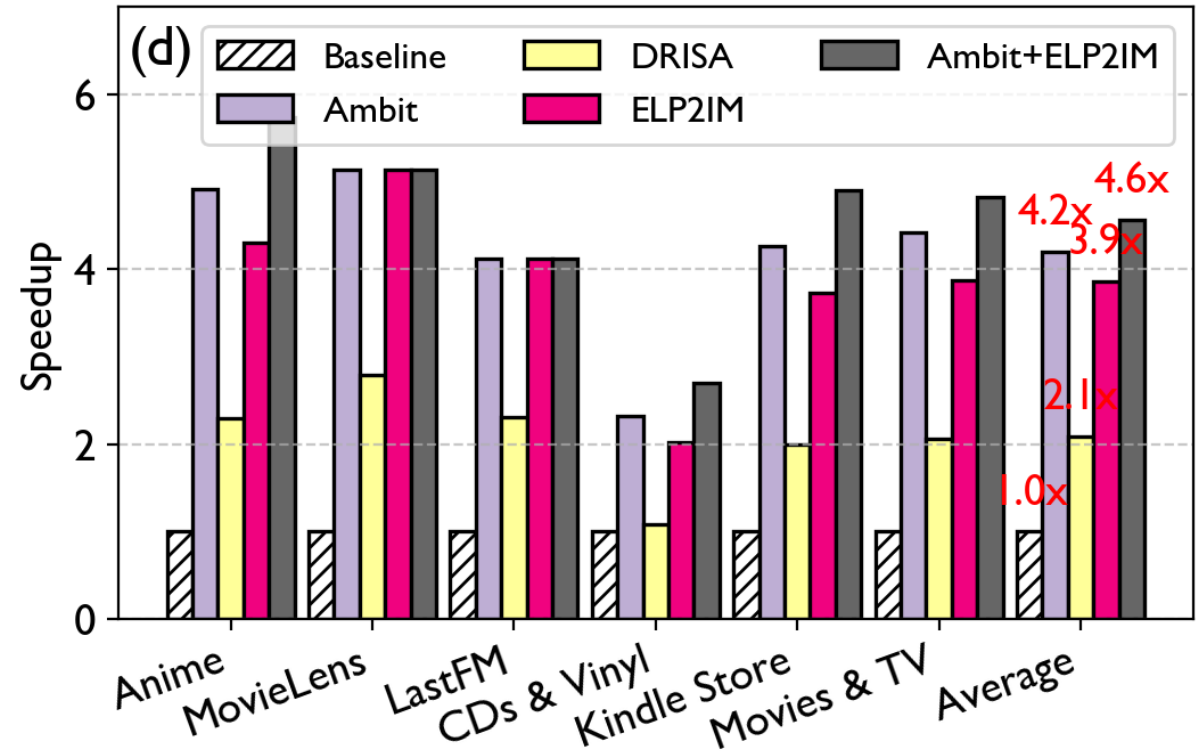
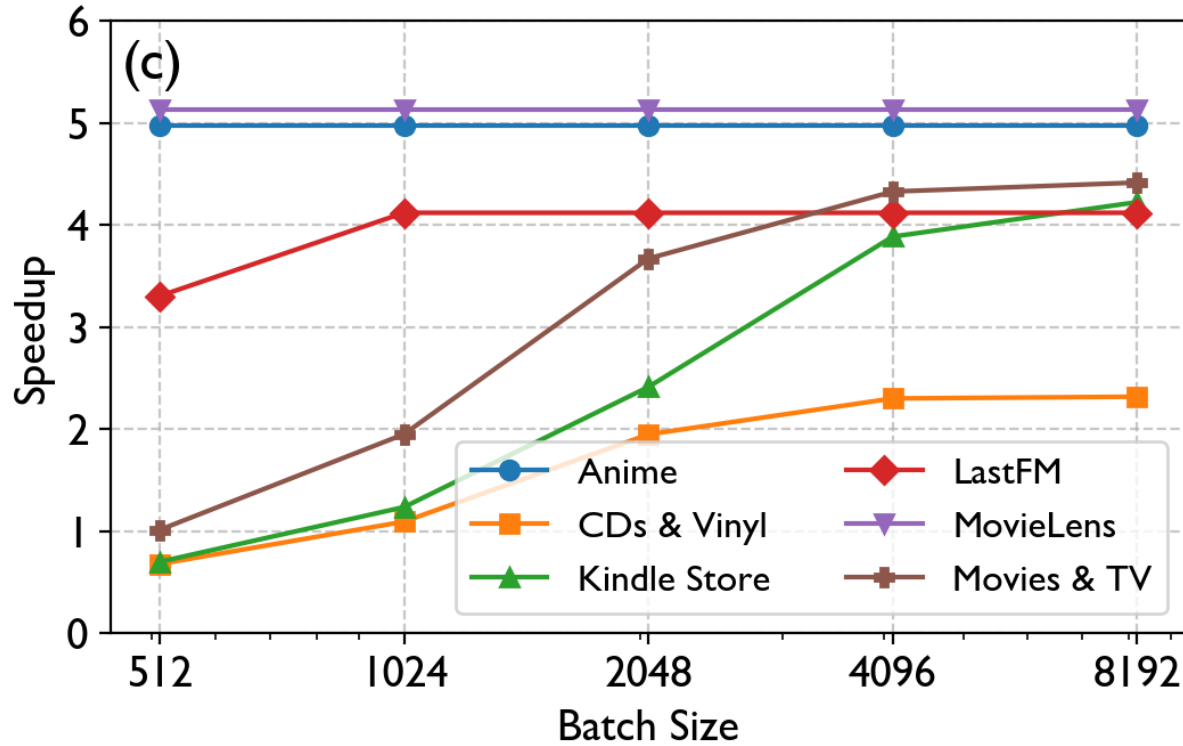
1. Eliminates all logical operations that do not reduce the number of rows.
2. Performs accumulation without requiring carry shifts.
3. Guarantees computational correctness through ECC-based reliability support.

Mapping	Carry-shift free	Read efficiency
Horizontal		✓
Vertical	✓	
PIM-SUM	✓	✓

# Performance--PIMSUM

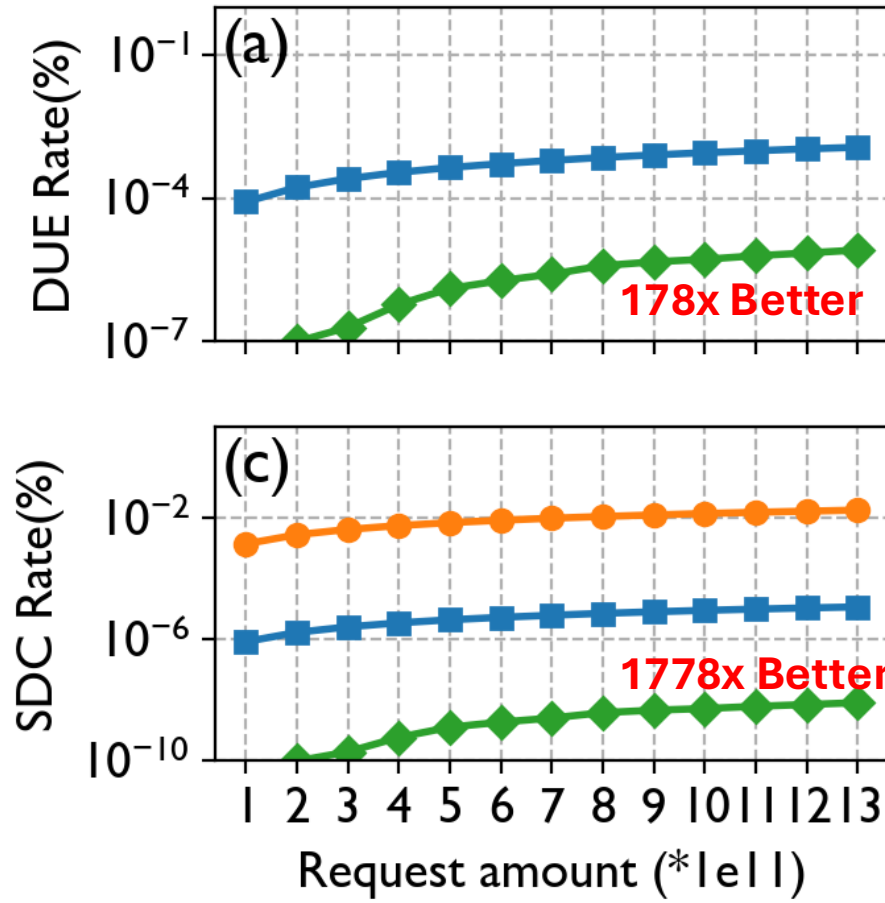


# Performance--PIMSUM

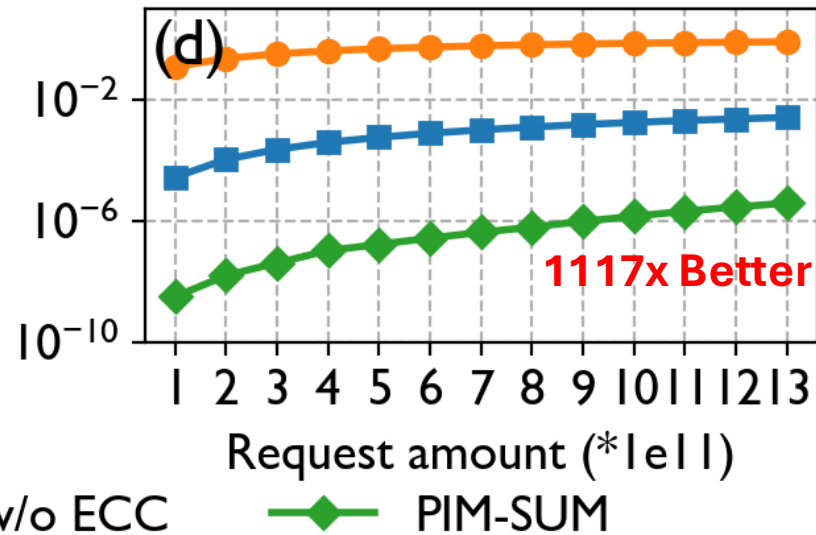
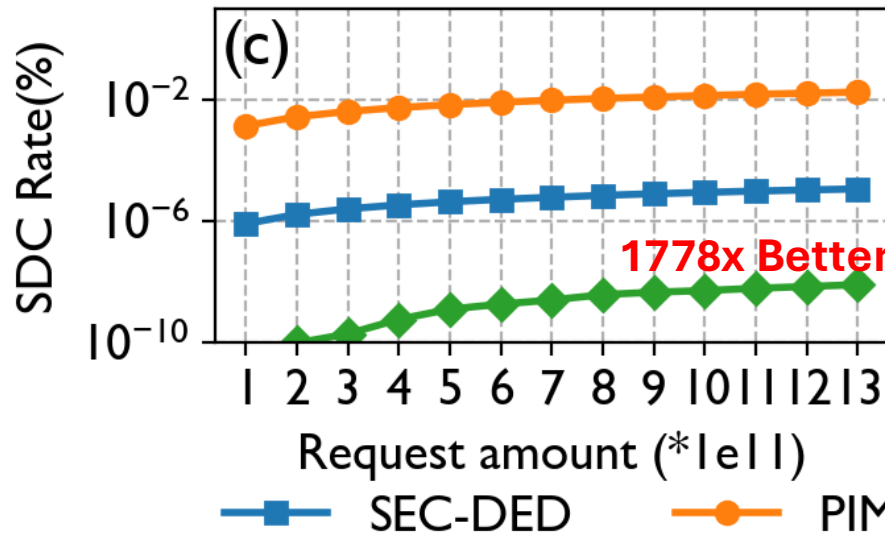
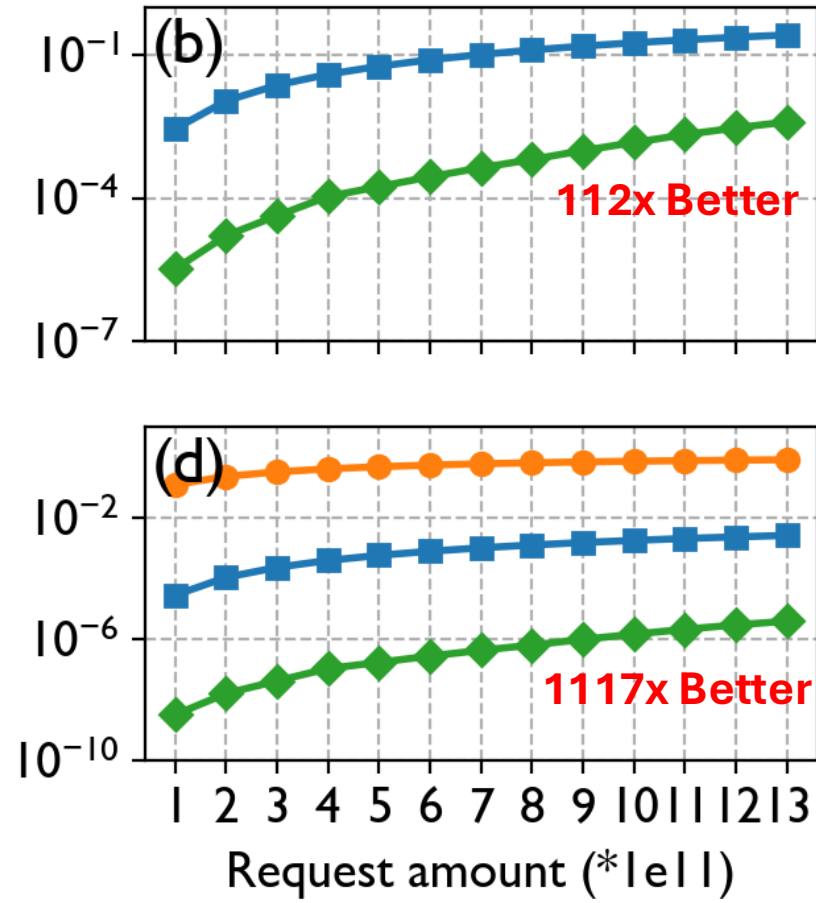


# Reliability--PIMSUM

Bit Flip Rate: 1e-6



Bit Flip Rate: 1e-8



■ SEC-DED    
 ● PIM w/o ECC    
 ◆ PIM-SUM

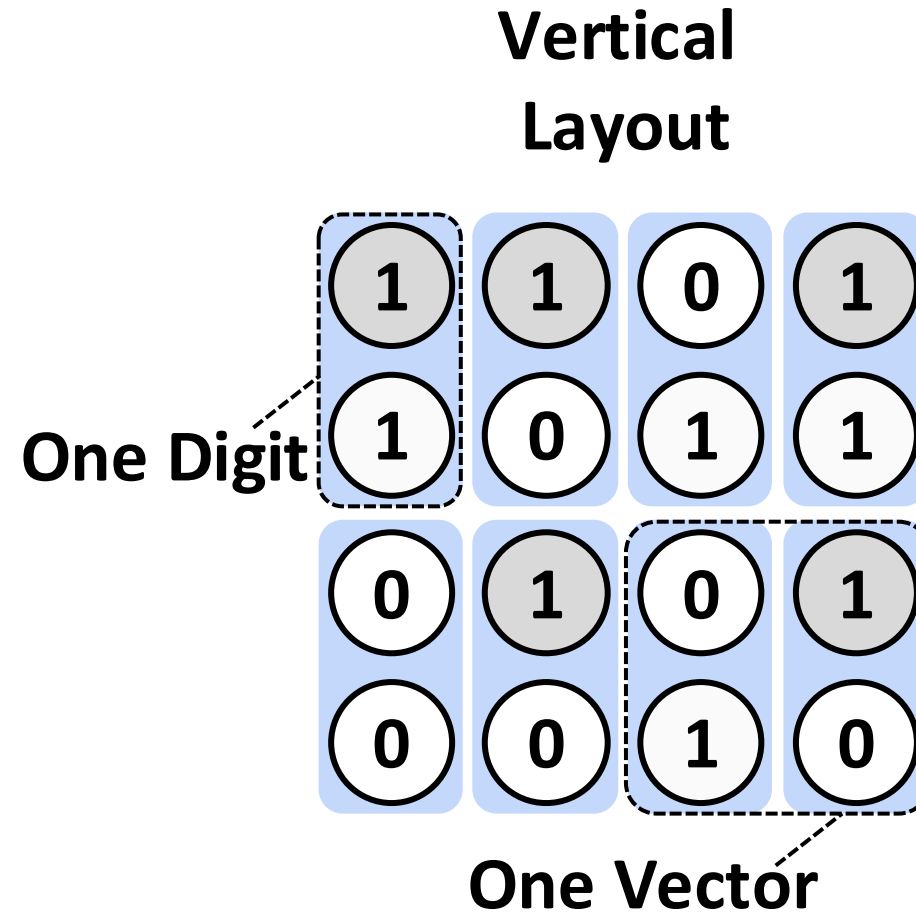
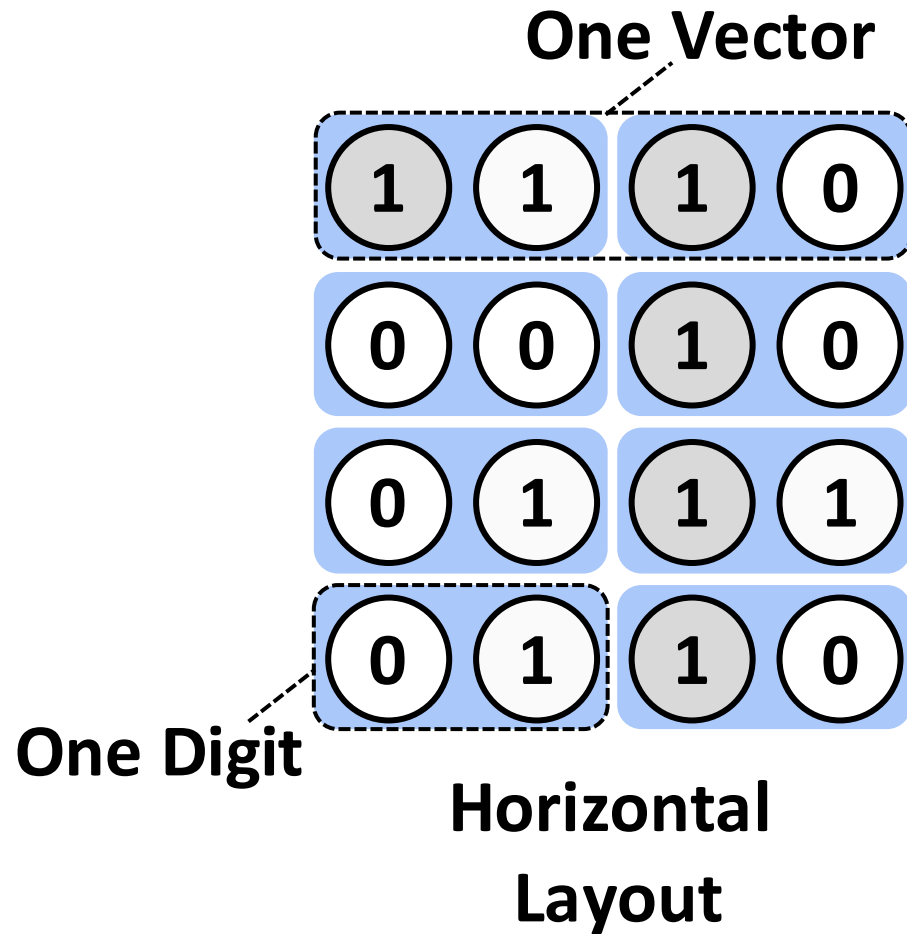
Thank you for listening  
Q&A

# Summary--Horizontal Layout

- Digit width: **w(=2)**
- Embedding vector element count: **n(=8)**
- Aggregation scale: **k(=4)**
- Cost: **w(carry shift) \* k (row summation) + 1 (mem\_read)**

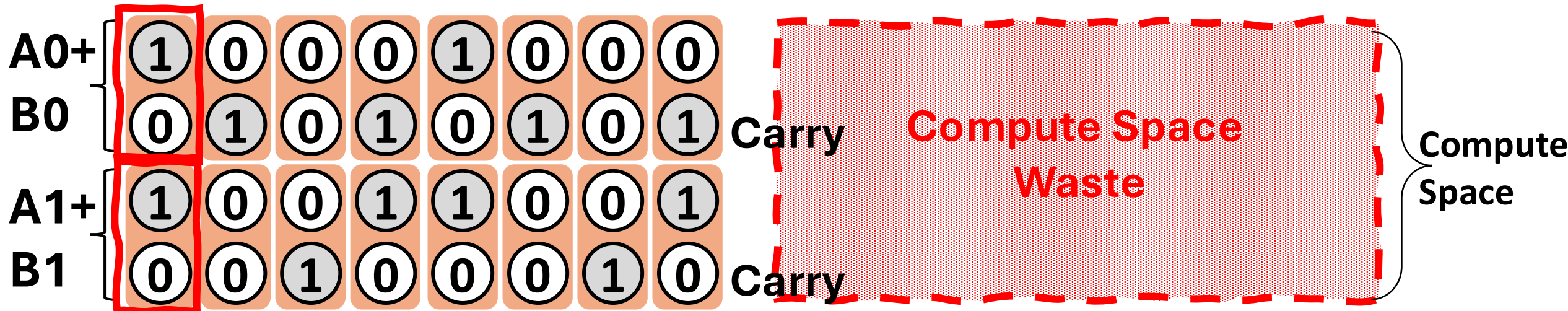
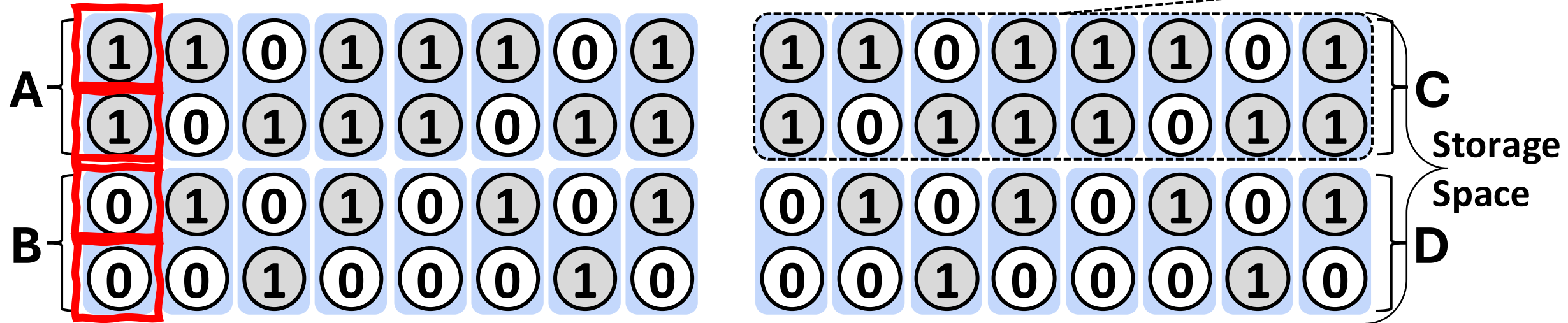


# Mapping



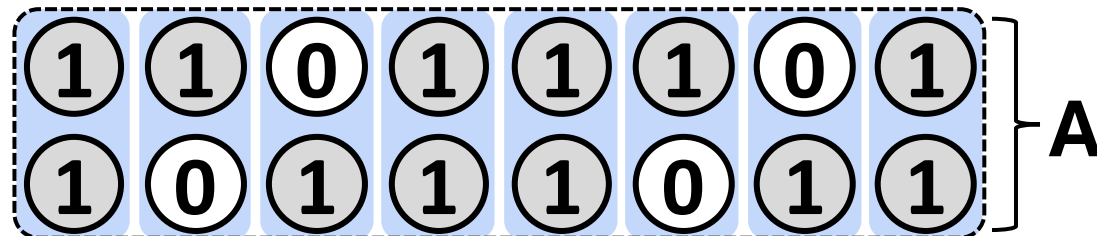
# Intuitive Approach--Vertical Layout

One Vector



# Summary--Vertical Layout

- Digit width: **w(=2)**
- Embedding vector element count: **n(=8)**
- Aggregation scale: **k(=4)**
- Cost: **k(rows) \* n(row serial) + w (mem\_read)**



# Evaluation--Methodology

Processor	32 cores, 4.0GHz, non-inclusive L1: 16KB, 64B cacheline, 8-way asso. L2: 256KB, 64B cacheline, 8-way asso. LLC: 1MB/core, 64B cacheline, 16-way asso.
Memory Controller	Page management: open-page Scheduler: FR-FCFS Data buffer latency: 4
DRAM Memory	DDR5-3200, 4GB, $\times 4$ I/O width 2 channels, 2 ranks, 8 banks-groups, 4 banks/bank-group 8Kb local row buffer CL-nRCD-nRP: 24-24-24 nCCDS-nCCDL-nCCDLWR: 8-8-32
PIM Boards	Ambit, DRISA, ELP2M, Ambit + ELP2M
Dataset Workloads	<b>Anime, MovieLens, LastFM,</b> <b>Amazon: CDs&amp;Vinyl, Kindle Store, TV&amp;movie</b>