

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

Nastaran Hajinazar*

Geraldo F. Oliveira*

Sven Gregorio

Joao Ferreira

Nika Mansouri Ghiasi

Minesh Patel

Mohammed Alser

Saugata Ghose

Juan Gómez-Luna

Onur Mutlu

SAFARI

ETH zürich



Executive Summary

- **Motivation**: Processing-using-Memory (PuM) architectures can efficiently perform bulk bitwise computation
- **Problem**: Existing PuM architectures are not widely applicable
 - Support only a limited and specific set of operations
 - Lack the flexibility to support new operations
 - Require significant changes to the DRAM subarray
- **Goals**: Design a processing-using-DRAM framework that:
 - Efficiently implements complex operations
 - Provides the flexibility to support new desired operations
 - Minimally changes the DRAM architecture
- **SIMDRAM**: An end-to-end processing-using-DRAM framework that provides the programming interface, the ISA, and the hardware support for:
 1. Efficiently computing complex operations
 2. Providing the ability to implement arbitrary operations as required
 3. Using a massively-parallel in-DRAM SIMD substrate that requires minimal changes to DRAM
- **Key Results**: SIMDRAM provides:
 - 88x and 5.8x the throughput and 257x and 31x the energy efficiency of a baseline CPU and a high-end GPU, respectively, for 16 in-DRAM operations
 - 21x and 2.1x the performance of the CPU and GPU for seven real-world applications

Outline

1. Processing-using-DRAM

2. Background

3. SIMD RAM

- Processing-using-DRAM Substrate
- SIMD RAM Framework

4. System Integration

5. Evaluation

6. Conclusion

Outline

1. Processing-using-DRAM

2. Background

3. SIMD RAM

- Processing-using-DRAM Substrate
- SIMD RAM Framework

4. System Integration

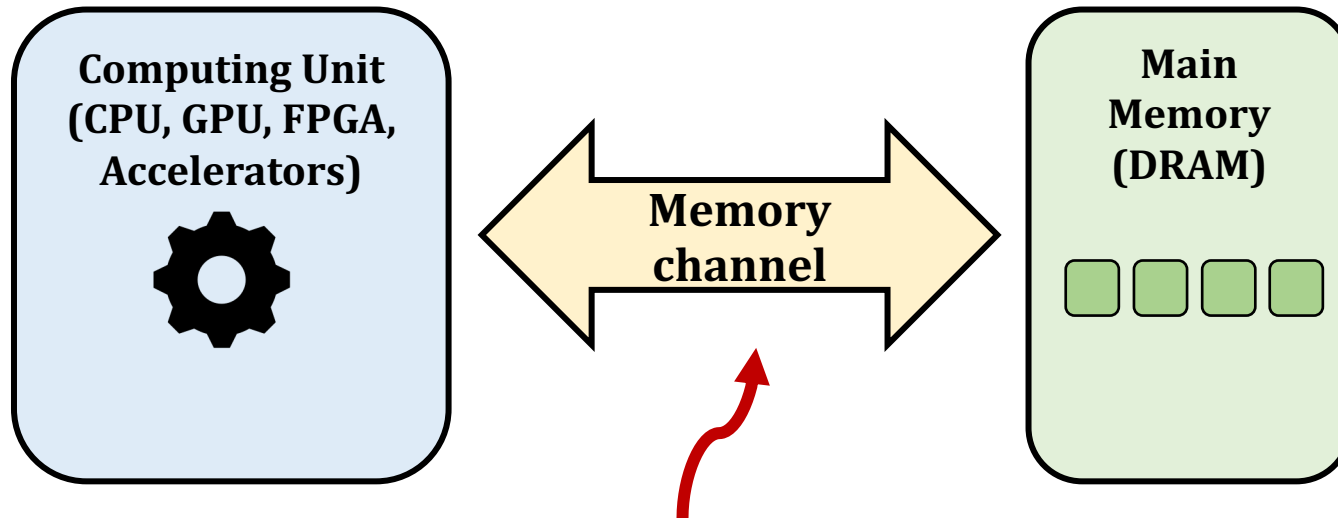
5. Evaluation

6. Conclusion

Data Movement Bottleneck

- Data movement is a major bottleneck

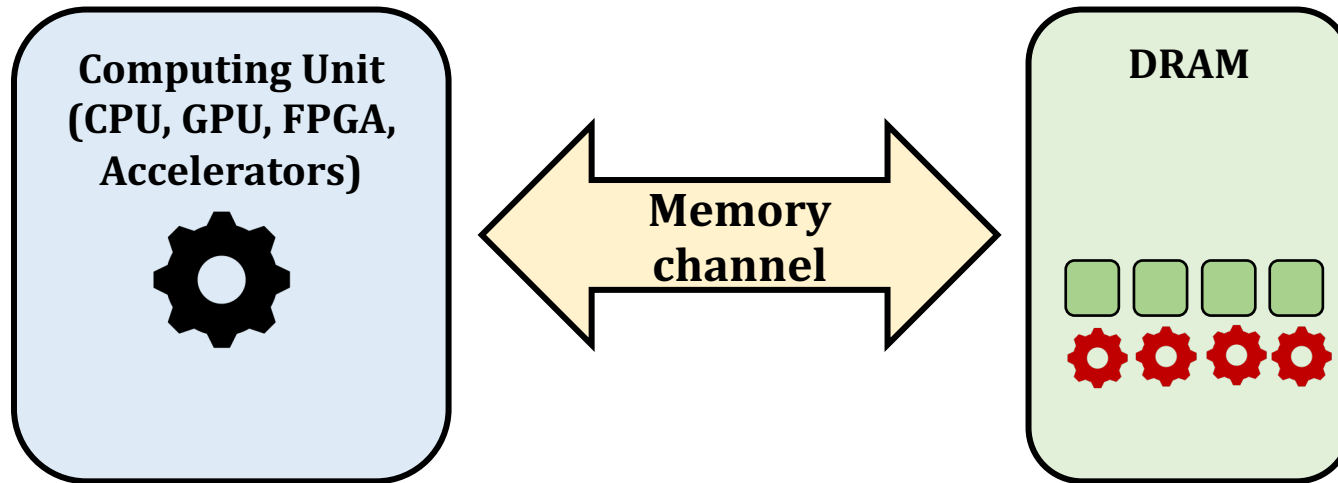
More than **60%** of the total system energy is spent on **data movement**¹



Bandwidth-limited and power-hungry memory channel

Processing-in-Memory (PIM)

- **Processing-in-Memory:** moves computation closer to where the data resides
 - **Reduces/eliminates** the need to move data between processor and DRAM



Processing-using-Memory (PuM)

- **PuM:** Exploits analog operation principles of the memory circuitry to perform computation
 - Leverages the **large internal bandwidth** and **parallelism** available inside the memory arrays
- A common approach for **PuM** architectures is to perform **bulk bitwise operations**
 - Simple logical operations (e.g., AND, OR, XOR)
 - More complex operations (e.g., addition, multiplication)

Outline

1. Processing-using-DRAM

2. Background

3. SIMDGRAM

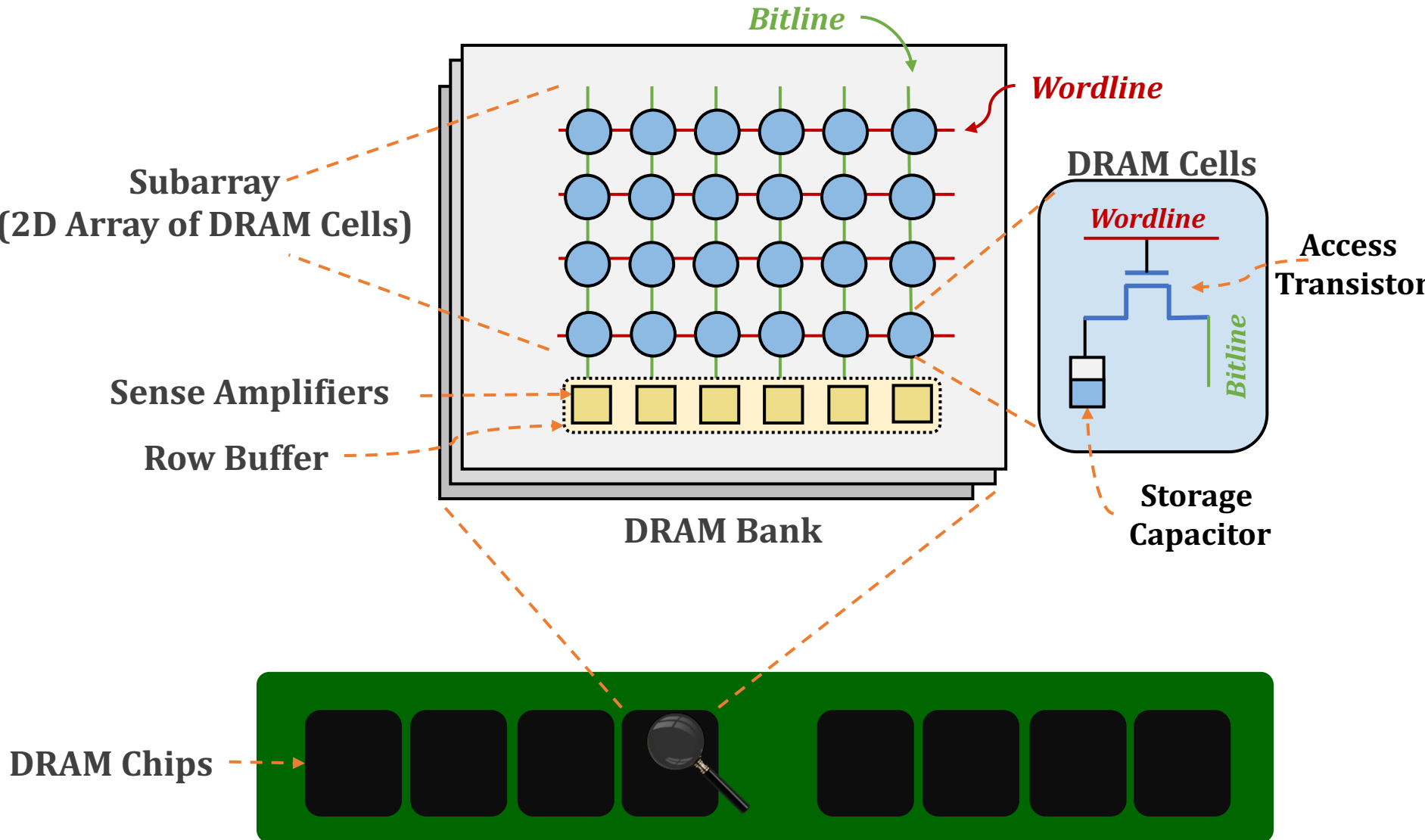
- Processing-using-DRAM Substrate
- SIMDGRAM Framework

4. System Integration

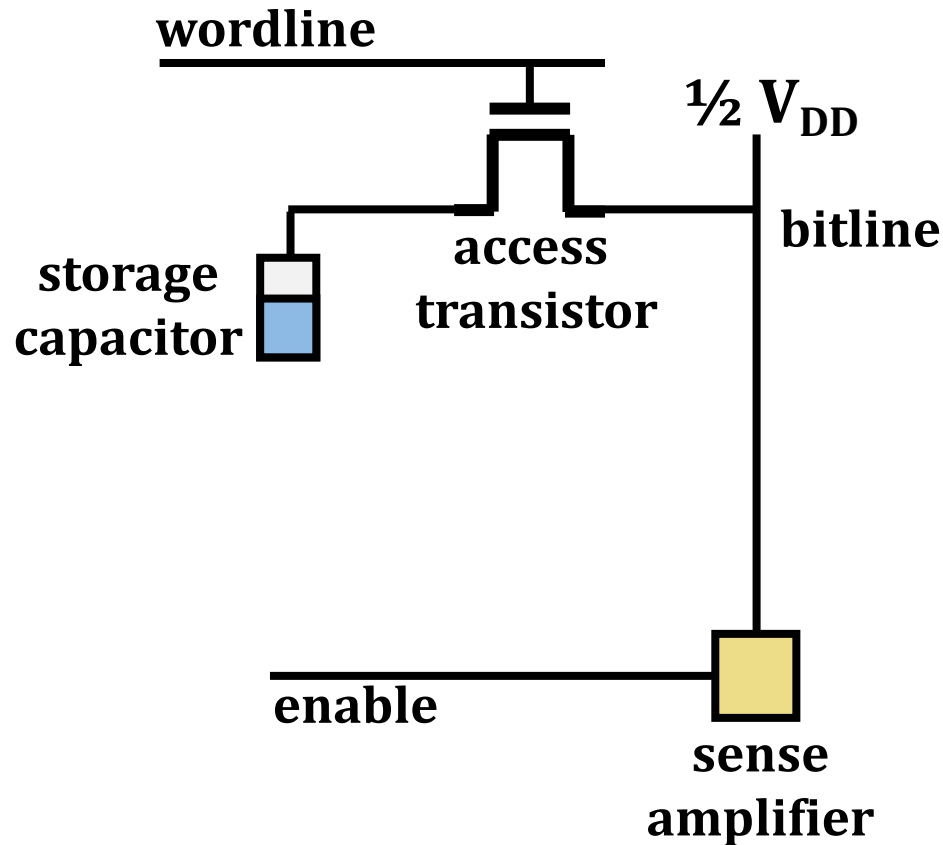
5. Evaluation

6. Conclusion

Inside a DRAM Chip



DRAM Cell Operation

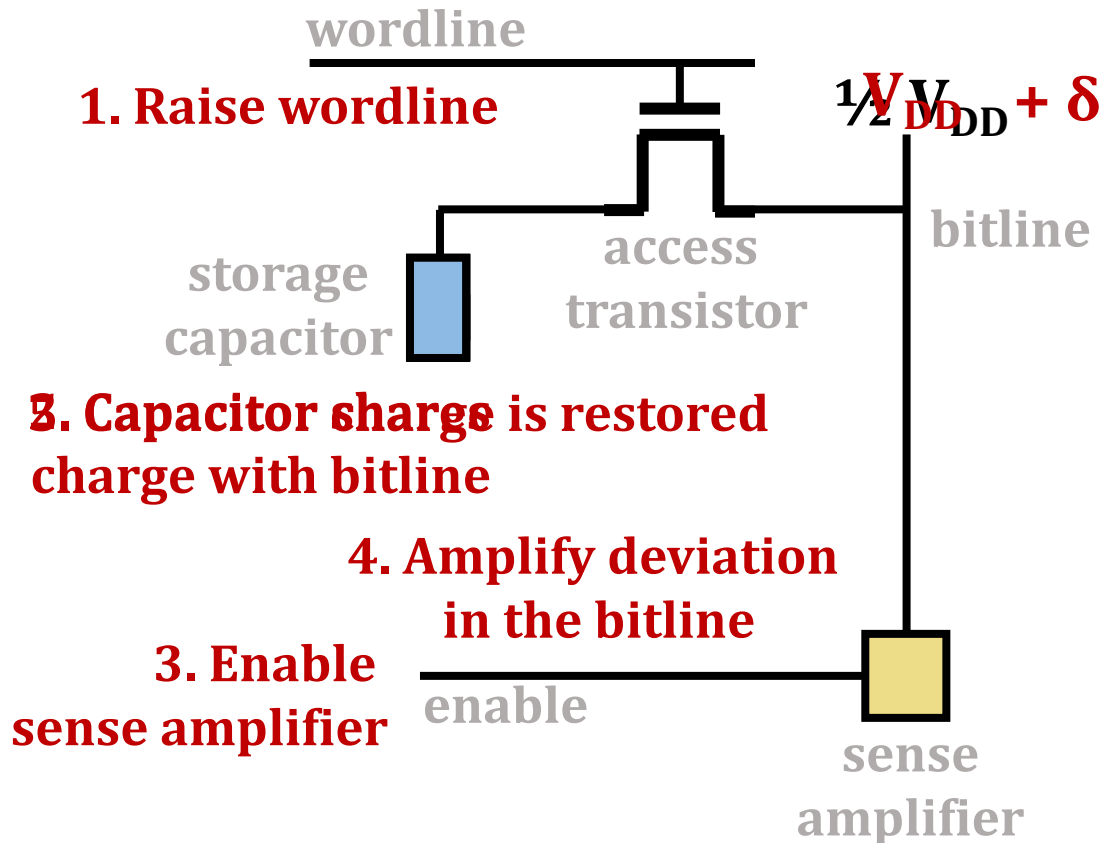


1. ACTIVATE (ACT)

2. READ/WRITE

3. PRECHARGE (PRE)

DRAM Cell Operation - ACTIVATE



1. Raise wordline

2. Capacitor charge is restored charge with bitline

4. Amplify deviation in the bitline

3. Enable sense amplifier

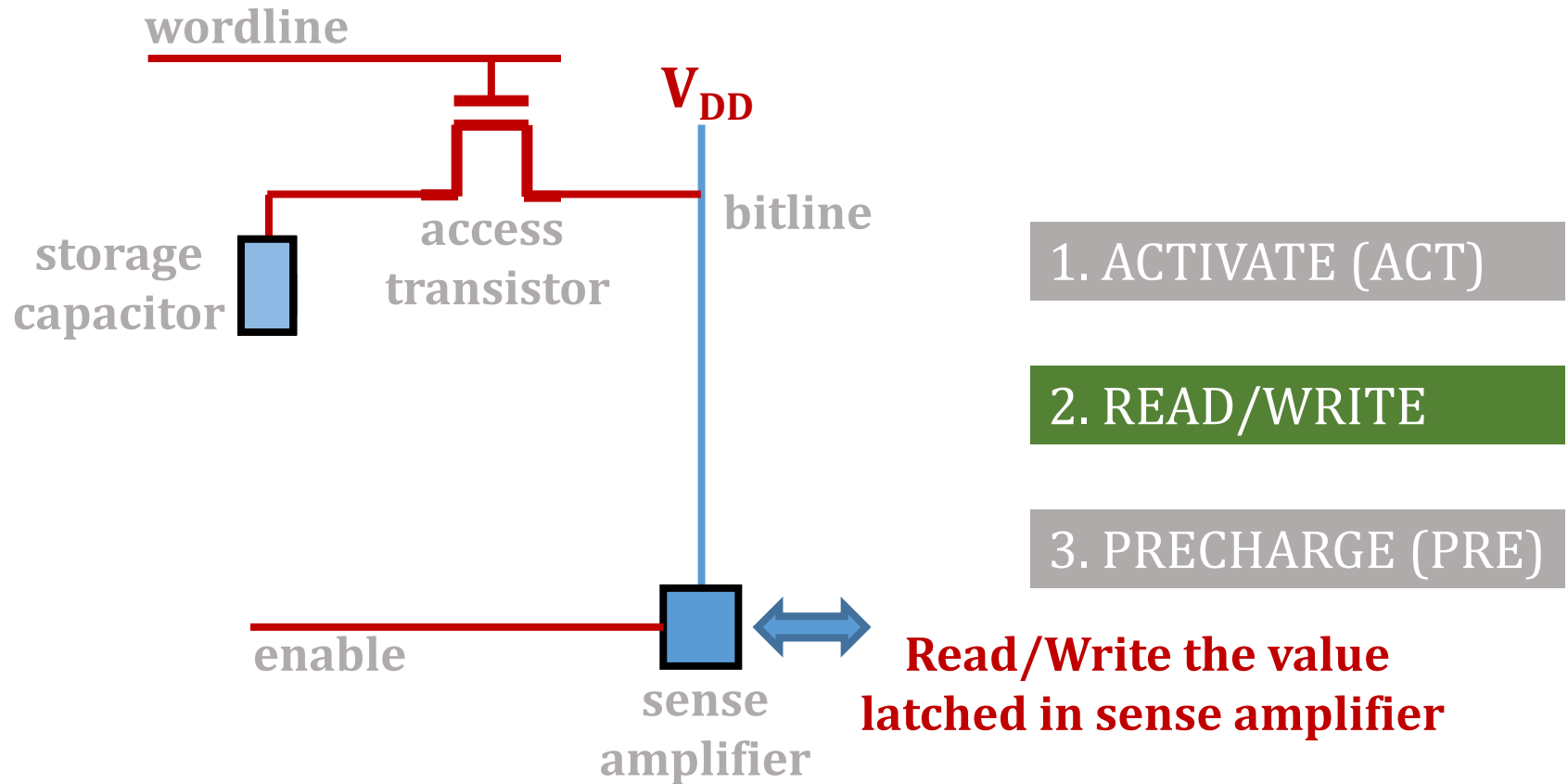
6. Row buffer stores the cell value

1. ACTIVATE (ACT)

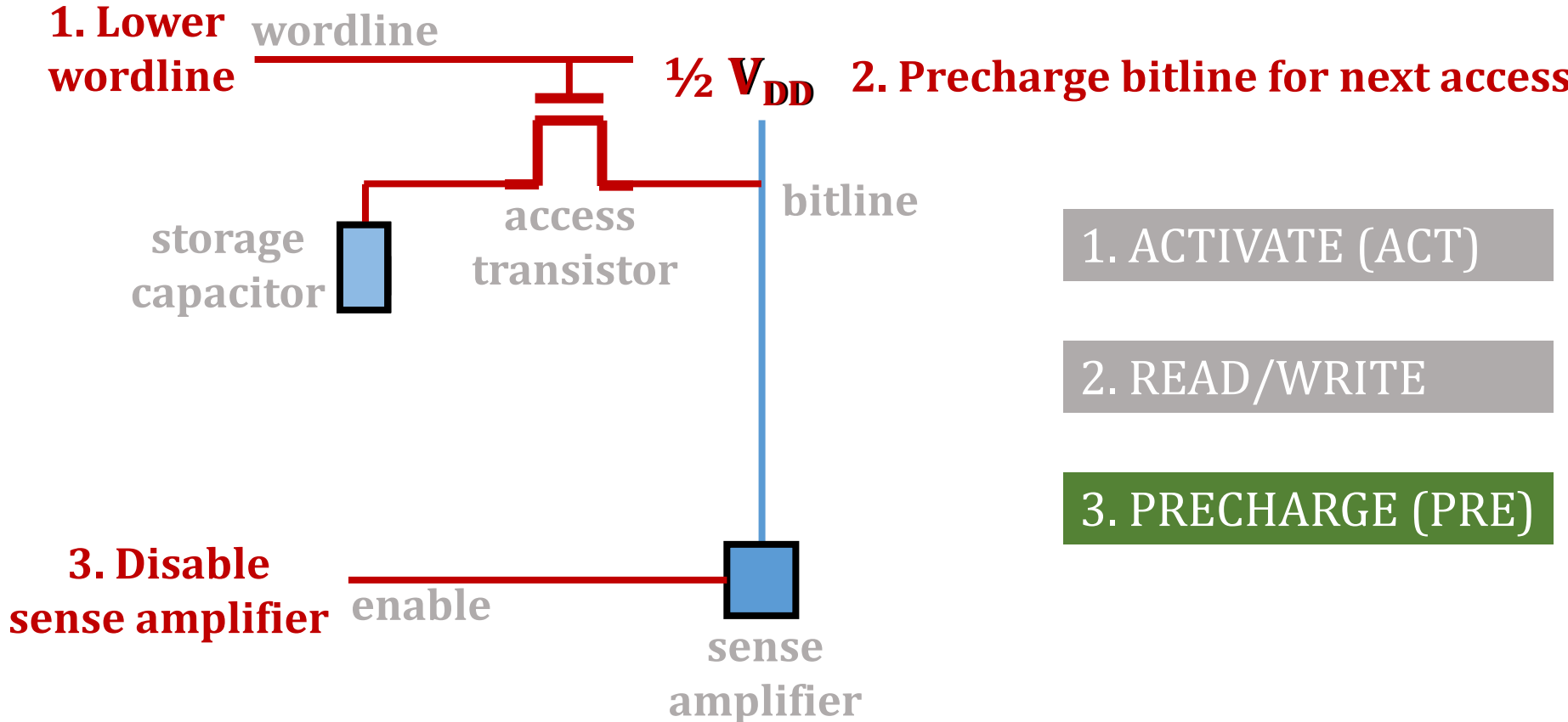
2. READ/WRITE

3. PRECHARGE (PRE)

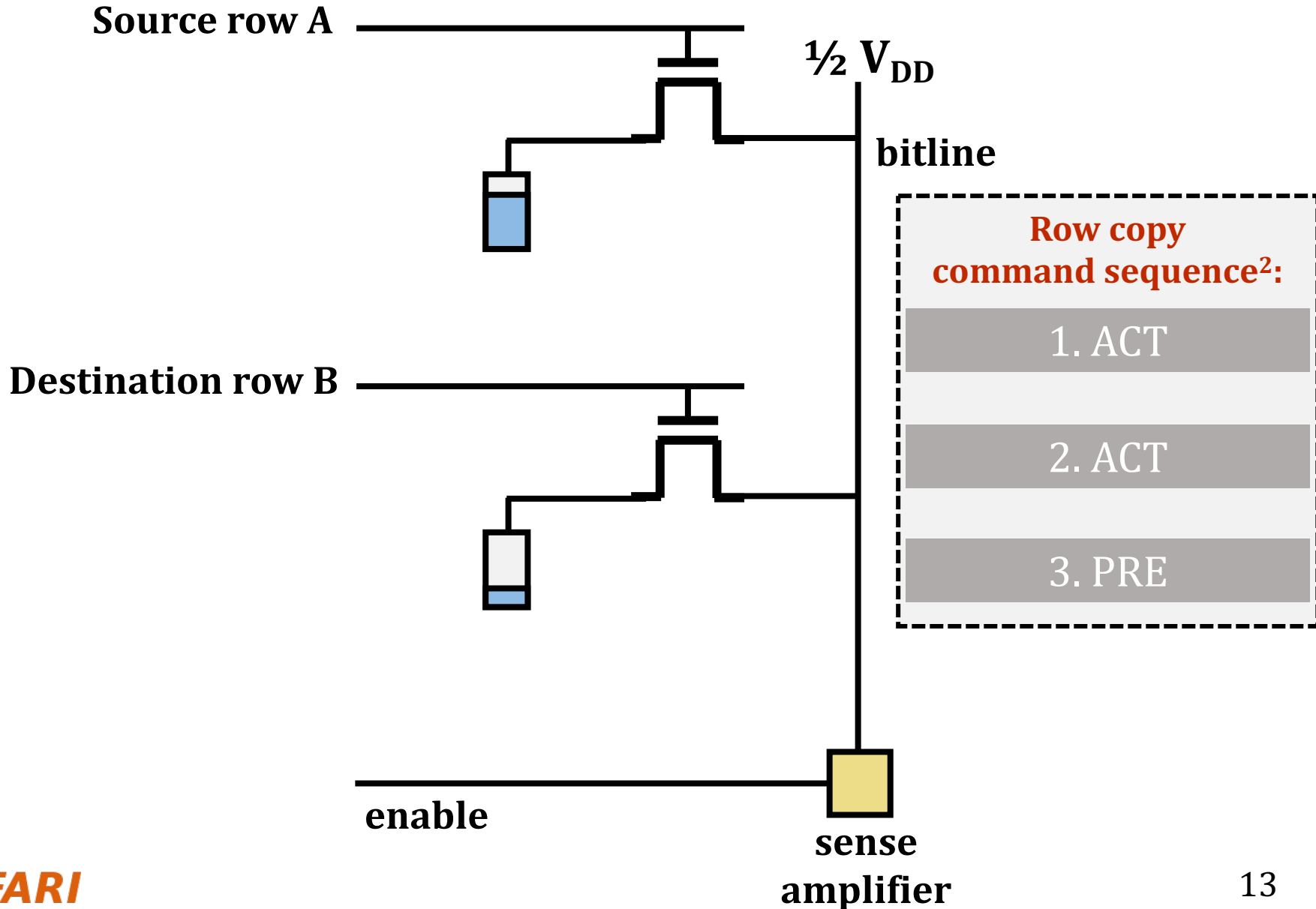
DRAM Cell Operation - READ/WRITE



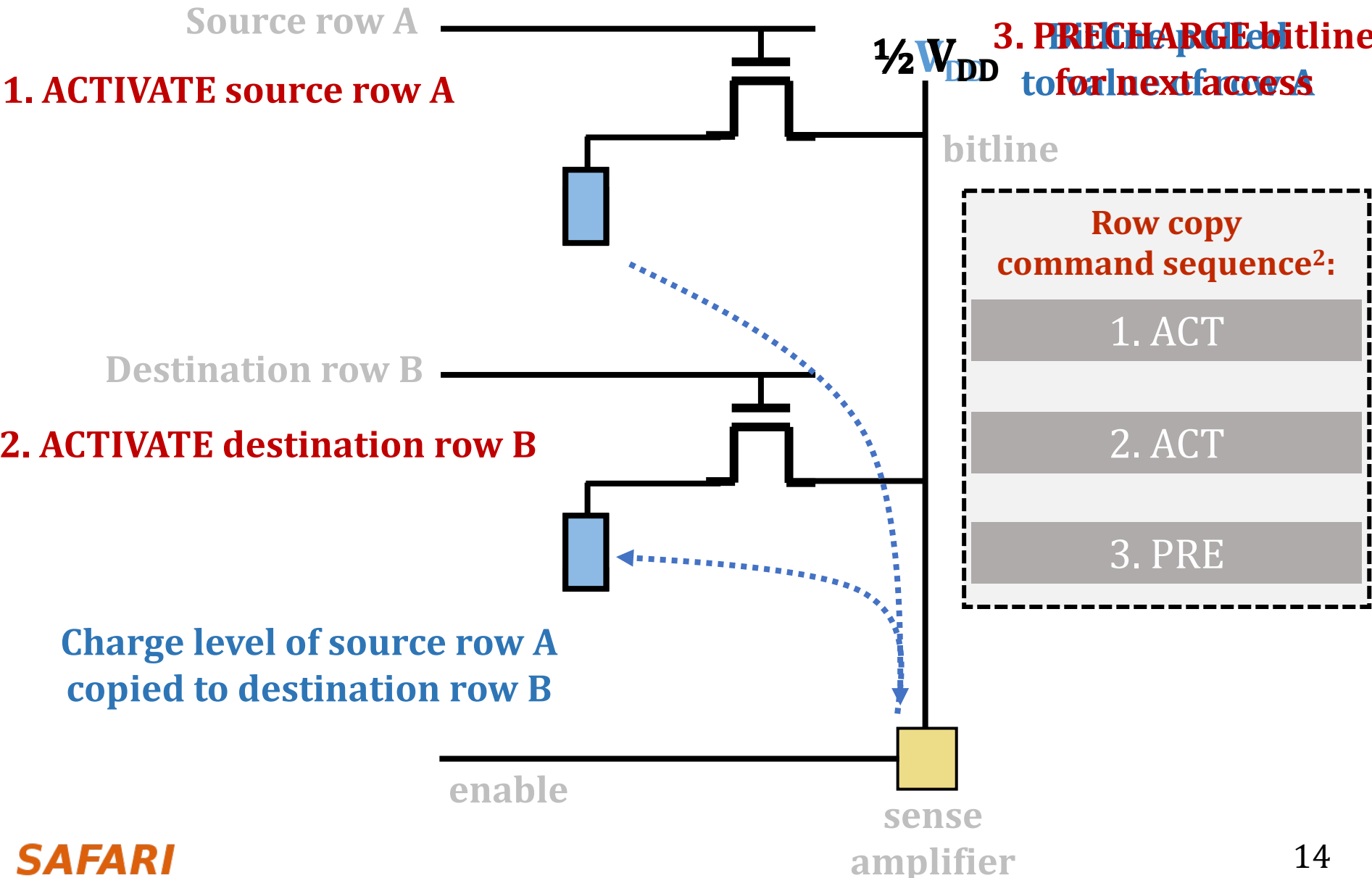
DRAM Cell Operation - PRECHARGE



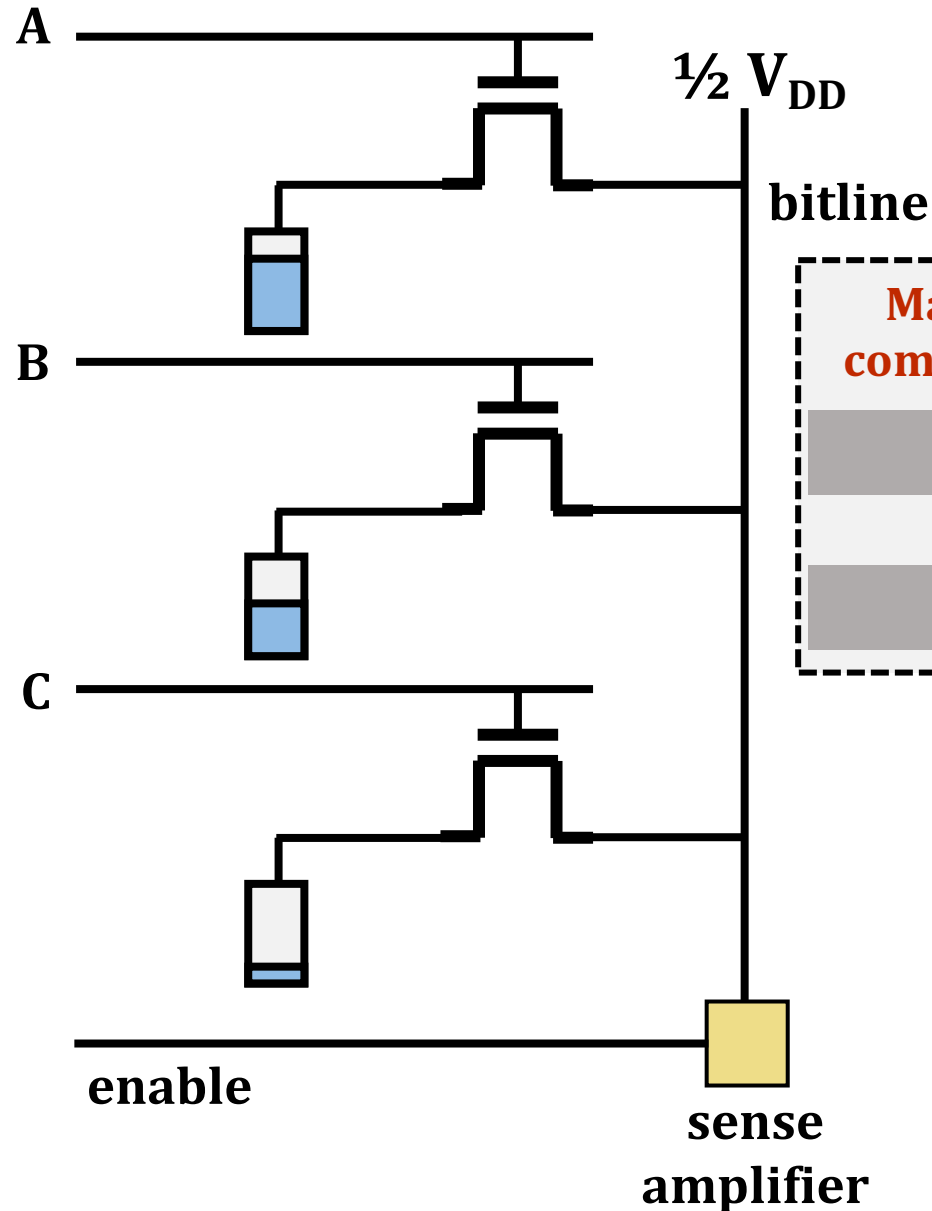
In-DRAM Row Copy



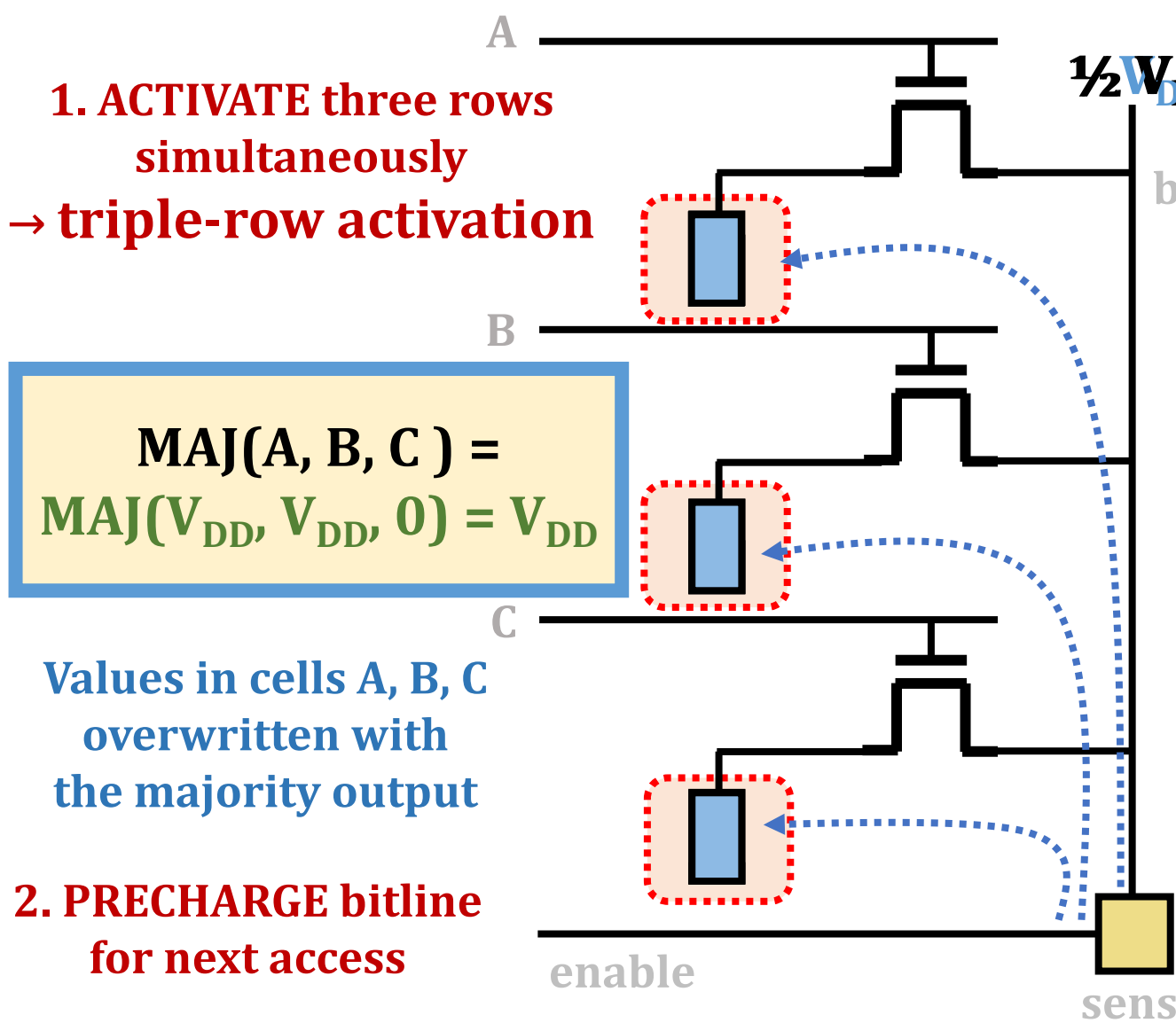
In-DRAM Row Copy: RowClone



Triple-Row Activation



Majority Function

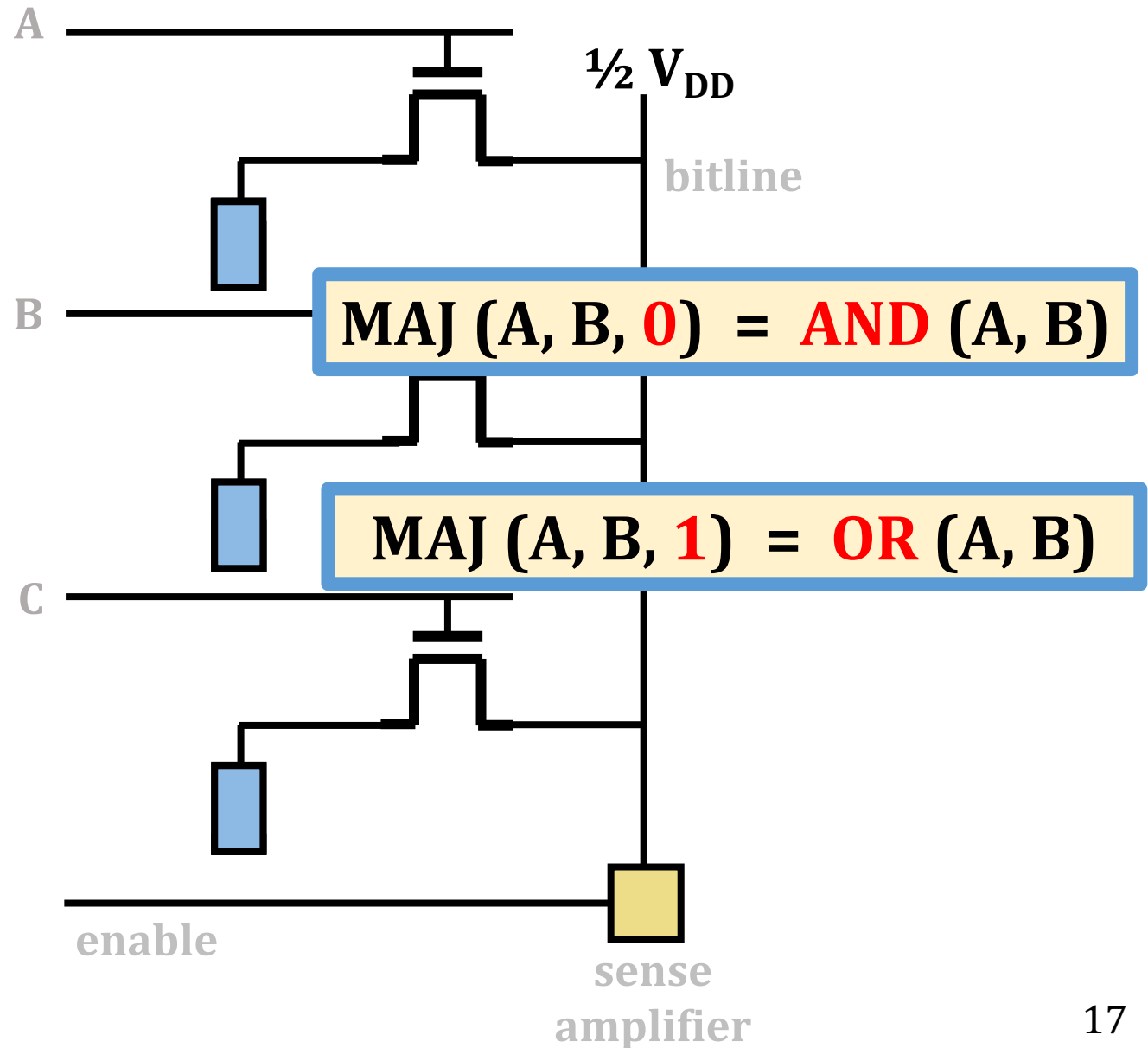


MAJ(A, B, C) =
 MAJ(V_{DD}, V_{DD}, 0) = V_{DD}

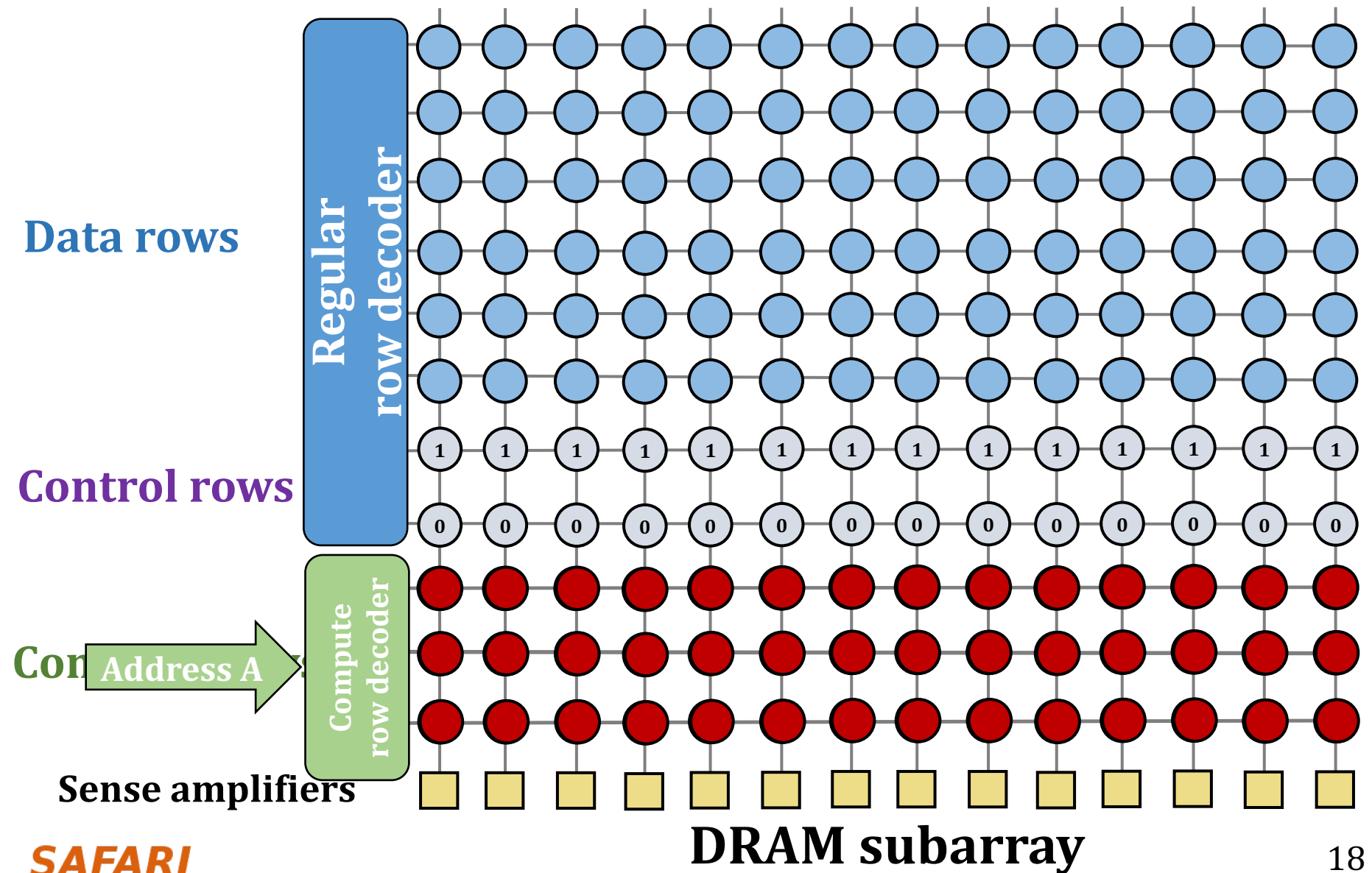
- Majority function command sequence³:
- 1. ACT
 - 2. PRE

³V. Seshadri et al., "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology", MICRO, 2017

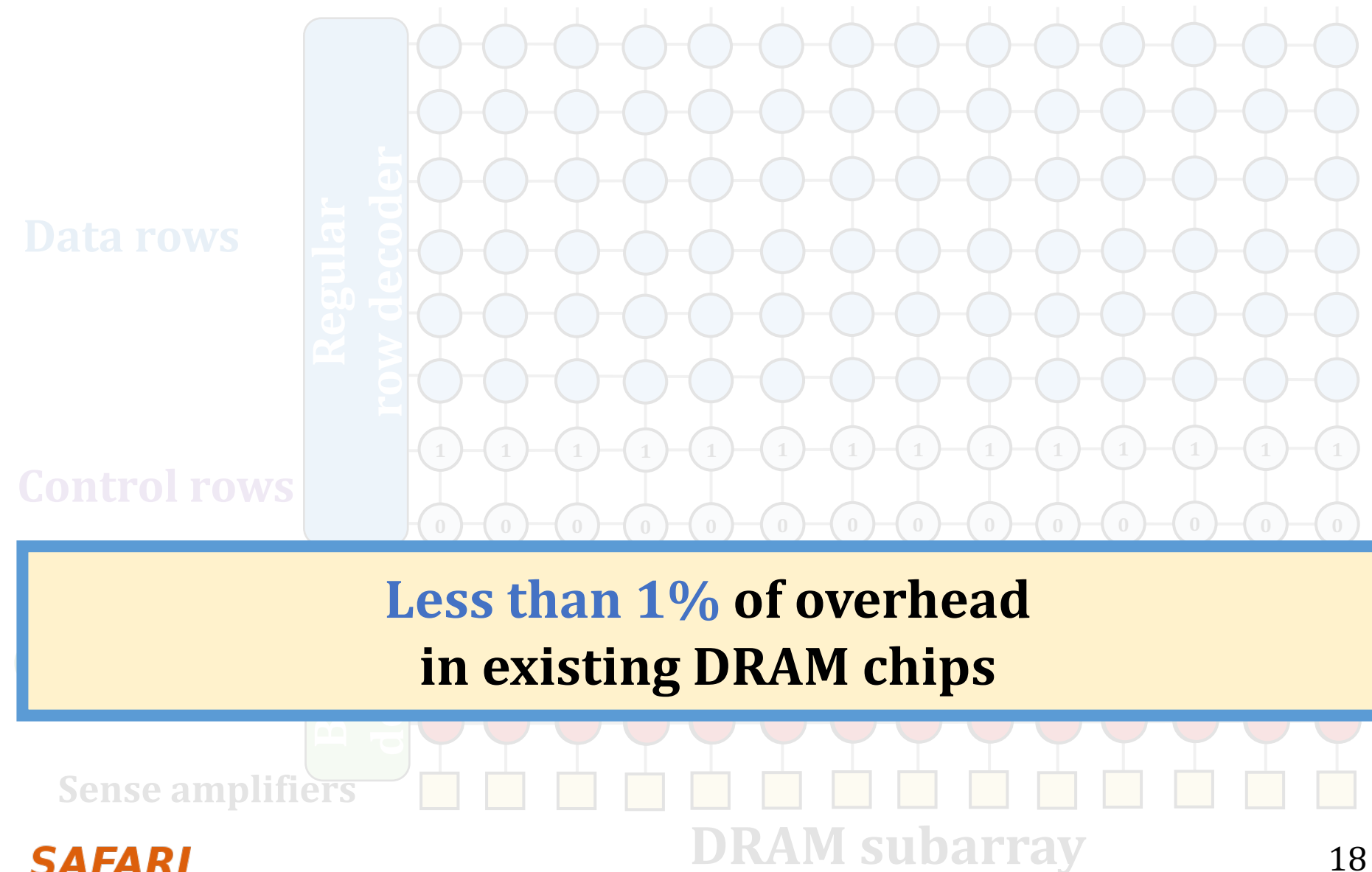
Ambit: In-DRAM Bulk Bitwise AND/OR



Ambit: Subarray Organization



Ambit: Subarray Organization



PuM: Prior Works

- DRAM and other memory technologies that are capable of performing **computation using memory**

Shortcomings:

- Support **only basic** operations (e.g., Boolean operations, addition)
 - Not widely applicable
- Support a **limited** set of operations
 - Lack the flexibility to support new operations
- Require **significant changes** to the DRAM
 - Costly (e.g., area, power)

PuM: Prior Works

- DRAM and other memory technologies that are capable of performing **computation using memory**

Shortcomings:

- Support **only basic** operations (e.g., Boolean operations, addition)

Need a framework that aids **general adoption of PuM, by:**

- **Efficiently implementing **complex operations****
- **Providing flexibility to support **new operations****

- Costly (e.g., area, power)

Our Goal

Goal: Design a PuM framework that

- **Efficiently** implements **complex** operations
- Provides the **flexibility** to support new desired operations
- **Minimally** changes the DRAM architecture

Outline

1. Processing-using-DRAM

2. Background

3. SIMD RAM

- Processing-using-DRAM Substrate
- SIMD RAM Framework

4. System Integration

5. Evaluation

6. Conclusion

Key Idea

- **SIMDRAM**: An end-to-end processing-using-DRAM framework that provides the **programming interface**, the **ISA**, and the **hardware support** for:
 - **Efficiently** computing **complex** operations in DRAM
 - Providing the ability to implement **arbitrary** operations as required
 - Using an **in-DRAM massively-parallel SIMD substrate** that requires **minimal** changes to DRAM architecture

Outline

1. Processing-using-DRAM

2. Background

3. SIMDGRAM

- Processing-using-DRAM Substrate
- SIMDGRAM Framework

4. System Integration

5. Evaluation

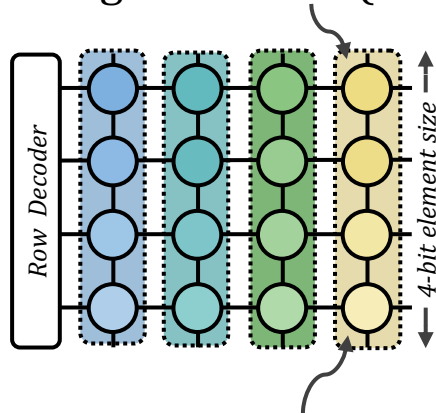
6. Conclusion

SIMDRAM: PuM Substrate

- SIMDRAM framework is built around a DRAM substrate that enables two techniques:

(1) Vertical data layout

most significant bit (MSB)



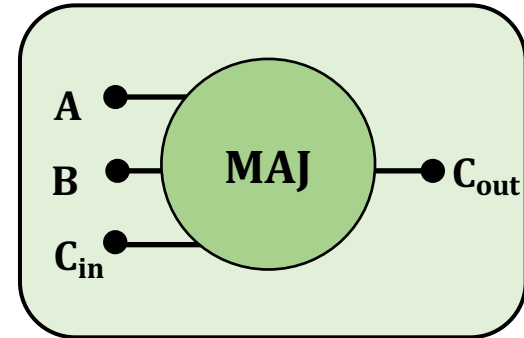
least significant bit (LSB)

Pros compared to the conventional **horizontal layout**:

- Implicit shift operation
- Massive parallelism

(2) Majority-based computation

$$C_{out} = AB + AC_{in} + BC_{in}$$



Pros compared to **AND/OR/NOT-based** computation:

- Higher performance
- Higher throughput
- Lower energy consumption

Outline

1. Processing-using-DRAM

2. Background

3. SIMDGRAM

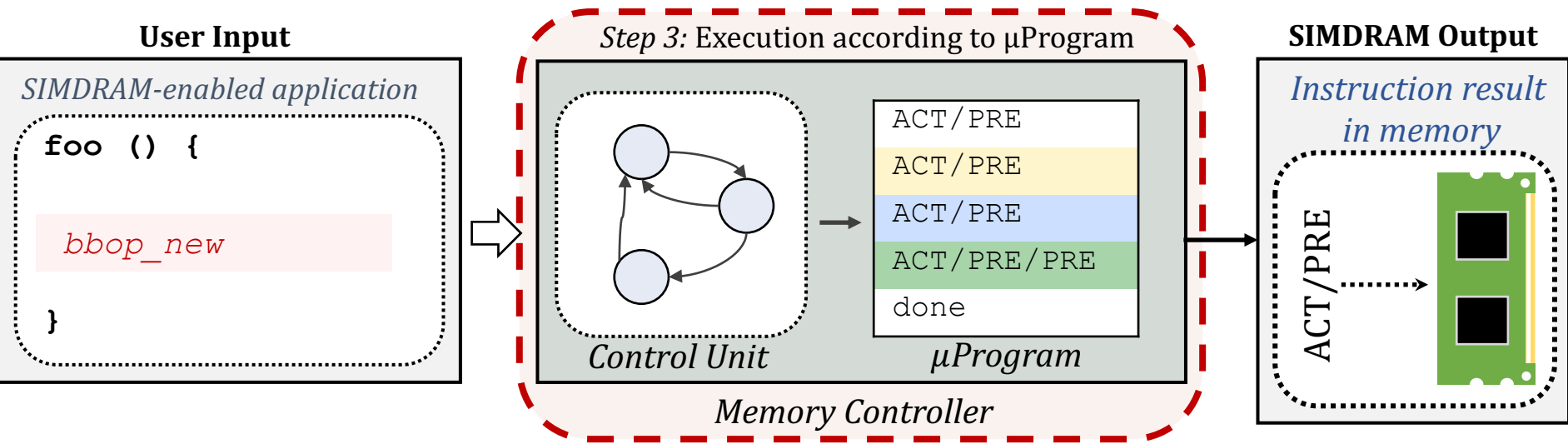
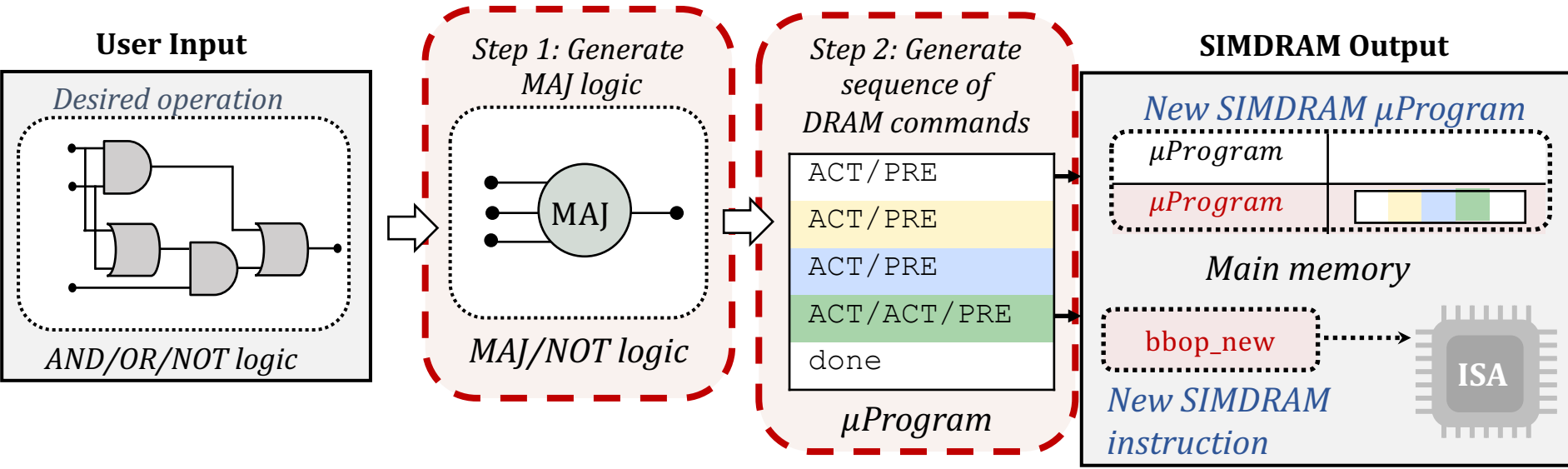
- Processing-using-DRAM Substrate
- SIMDGRAM Framework

4. System Integration

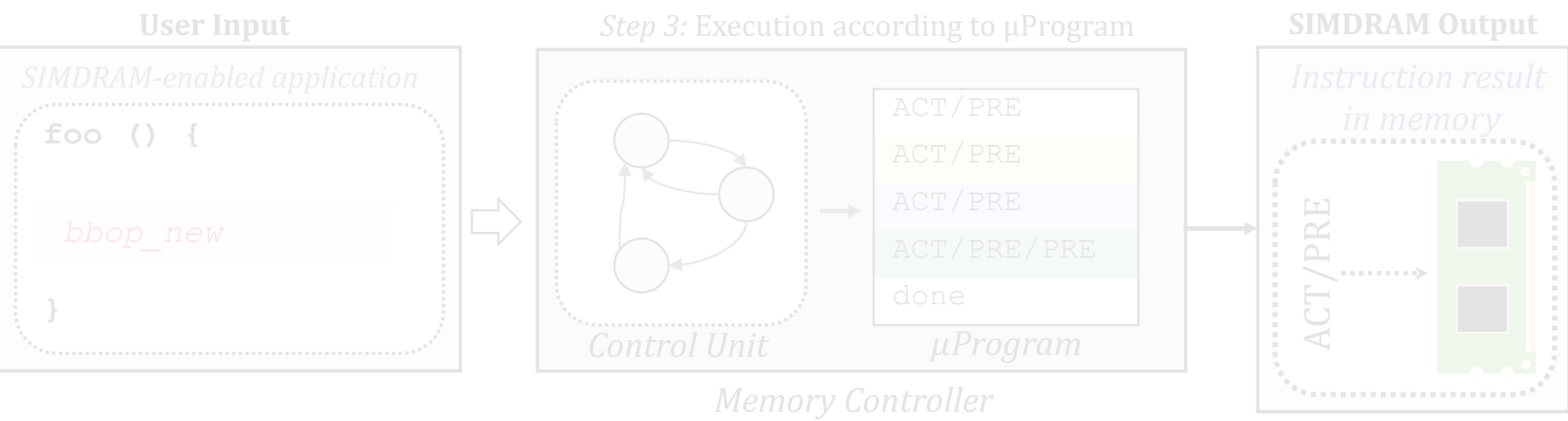
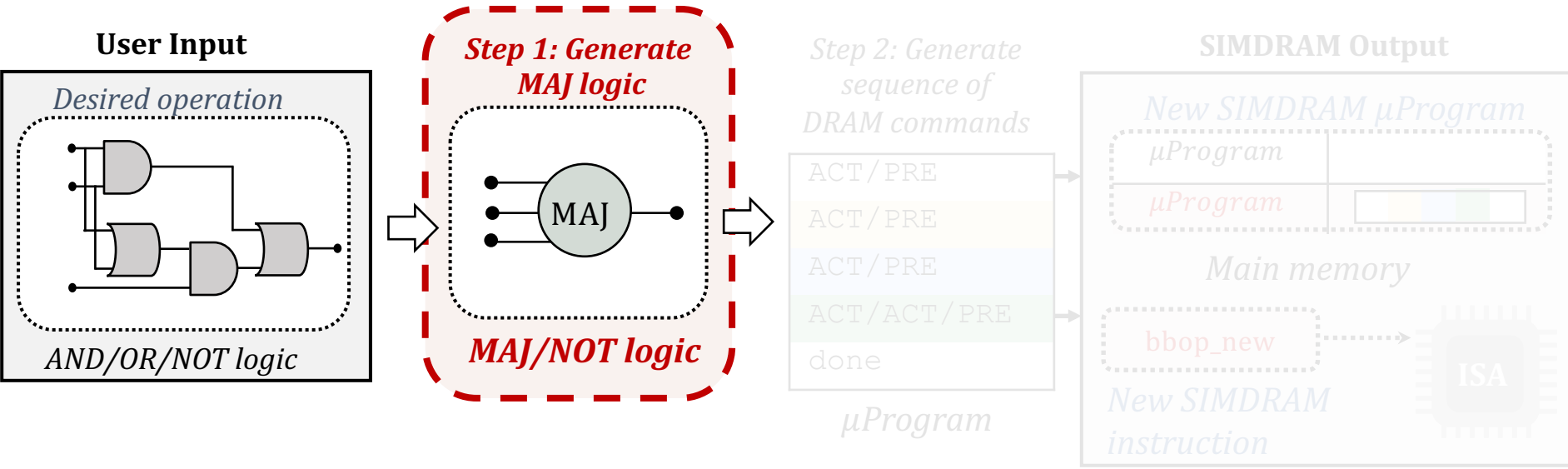
5. Evaluation

6. Conclusion

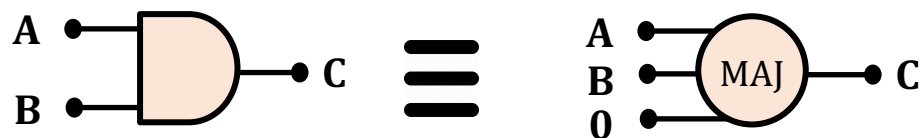
SIMDRAM Framework



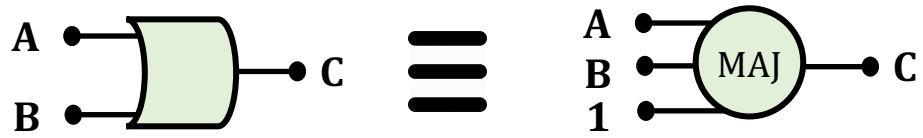
SIMDRAM Framework: Step 1



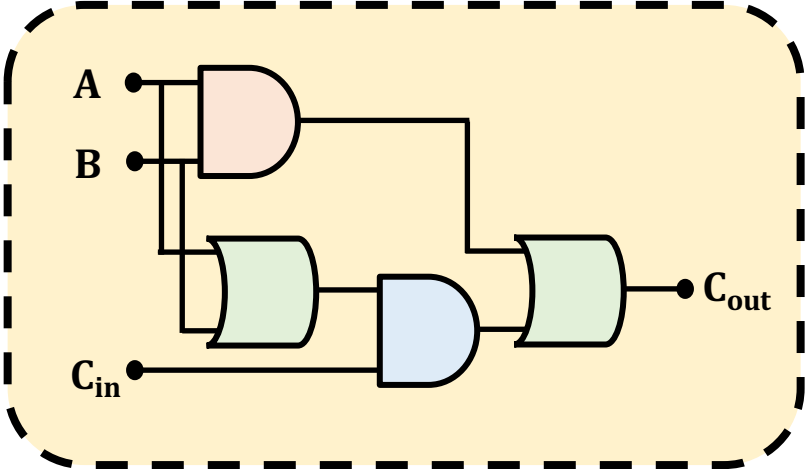
Step 1: Naïve MAJ/NOT Implementation



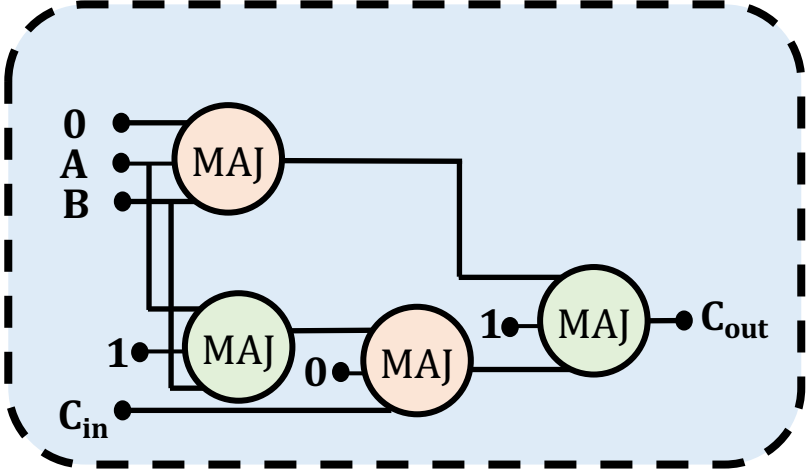
output is "1" only when A = B = "1"



output is "0" only when A = B = "0"

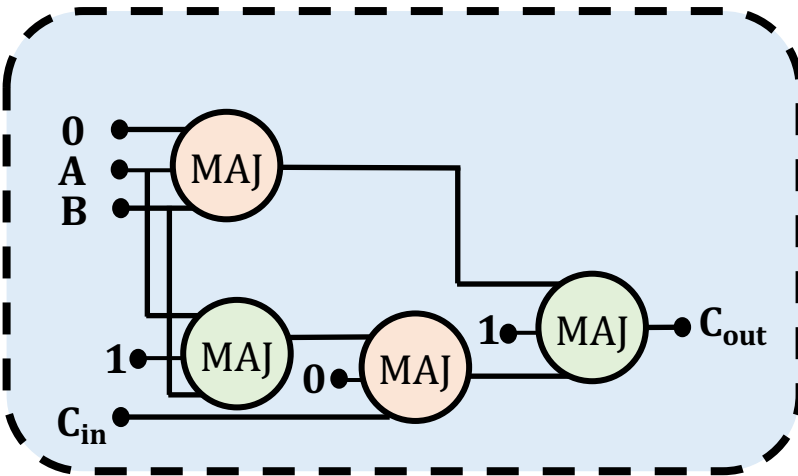


Part 1

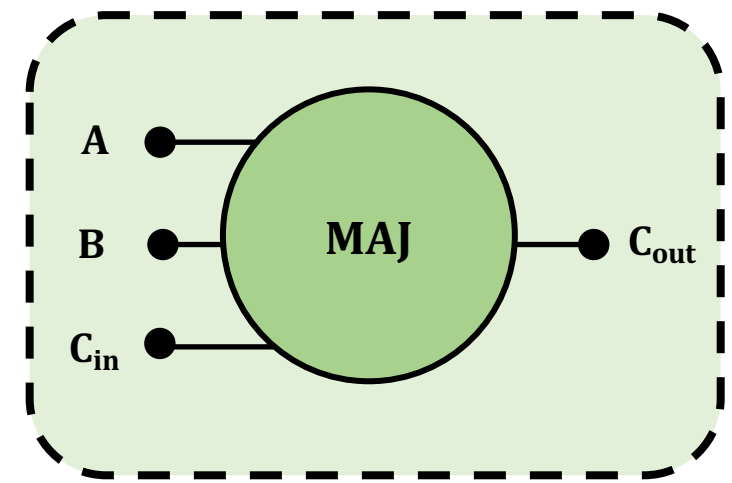
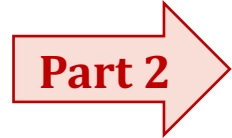


Naïvely converting AND/OR/NOT-implementation to MAJ/NOT-implementation leads to an **unoptimized circuit**

Step 1: Efficient MAJ/NOT Implementation



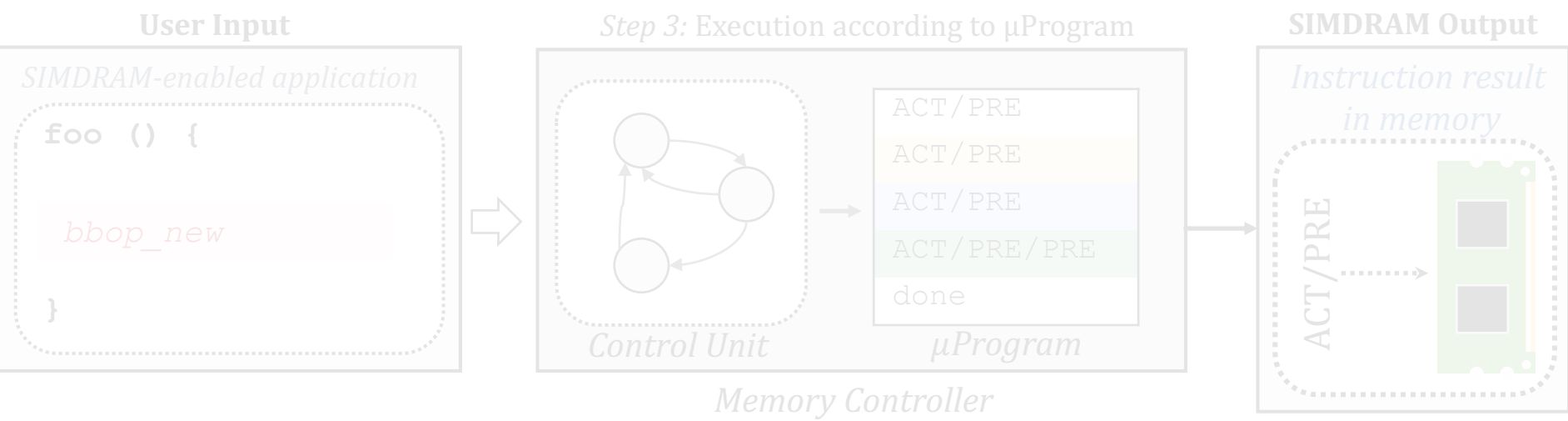
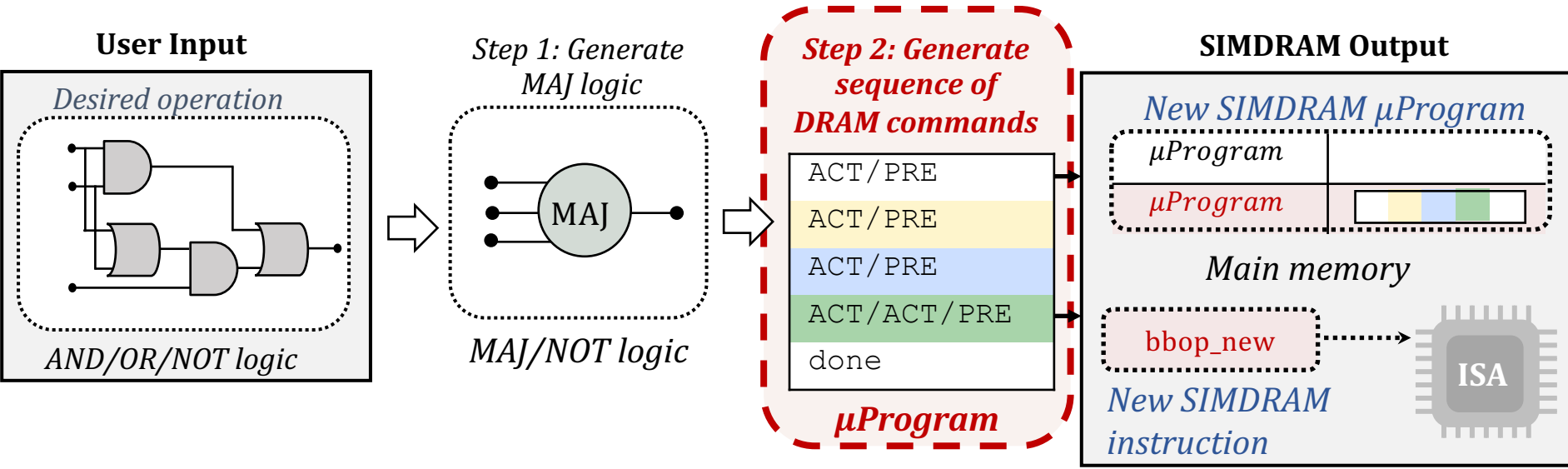
Greedy optimization algorithm⁴



Step 1 generates an **optimized MAJ/NOT-implementation** of the desired operation

⁴ L. Amarù et al, "Majority-Inverter Graph: A Novel Data-Structure and Algorithms for Efficient Logic Optimization", DAC, 2014.

SIMDRAM Framework: Step 2



Step 2: μ Program Generation

- **μ Program:** A series of **microarchitectural operations** (e.g., ACT/PRE) that SIMD RAM uses to execute **SIMDRAM operation in DRAM**
- **Goal of Step 2:** To generate the **μ Program** that **executes** the desired SIMD RAM operation **in DRAM**

Task 1: Allocate DRAM rows to the operands

Task 2: Generate μ Program

Step 2: μ Program Generation

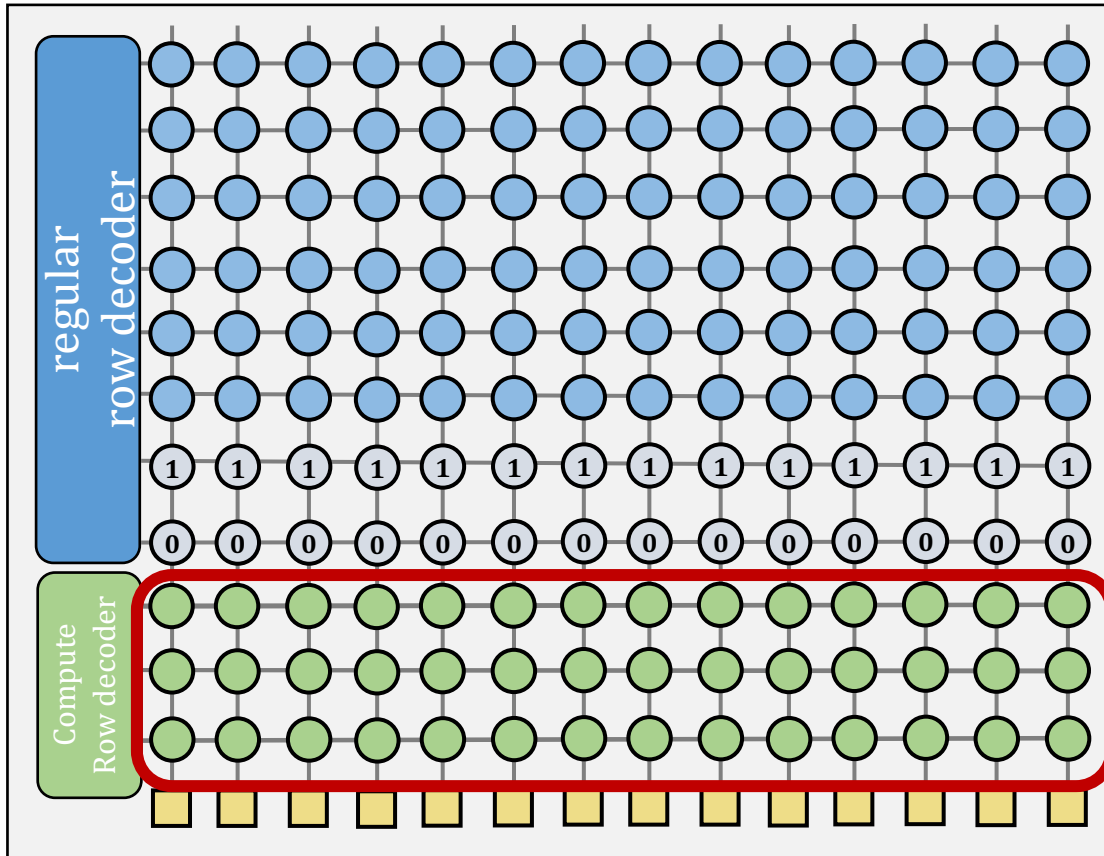
- **μ Program:** A series of **microarchitectural operations** (e.g., ACT/PRE) that SIMD RAM uses to execute **SIMDRAM operation in DRAM**
- **Goal of Step 2:** To generate the **μ Program** that executes the desired SIMD RAM operation in DRAM

Task 1: Allocate DRAM rows to the operands

Task 2: Generate μ Program

Task 1: Allocating DRAM Rows to Operands

- Allocation algorithm considers **two constraints** specific to processing-using-DRAM



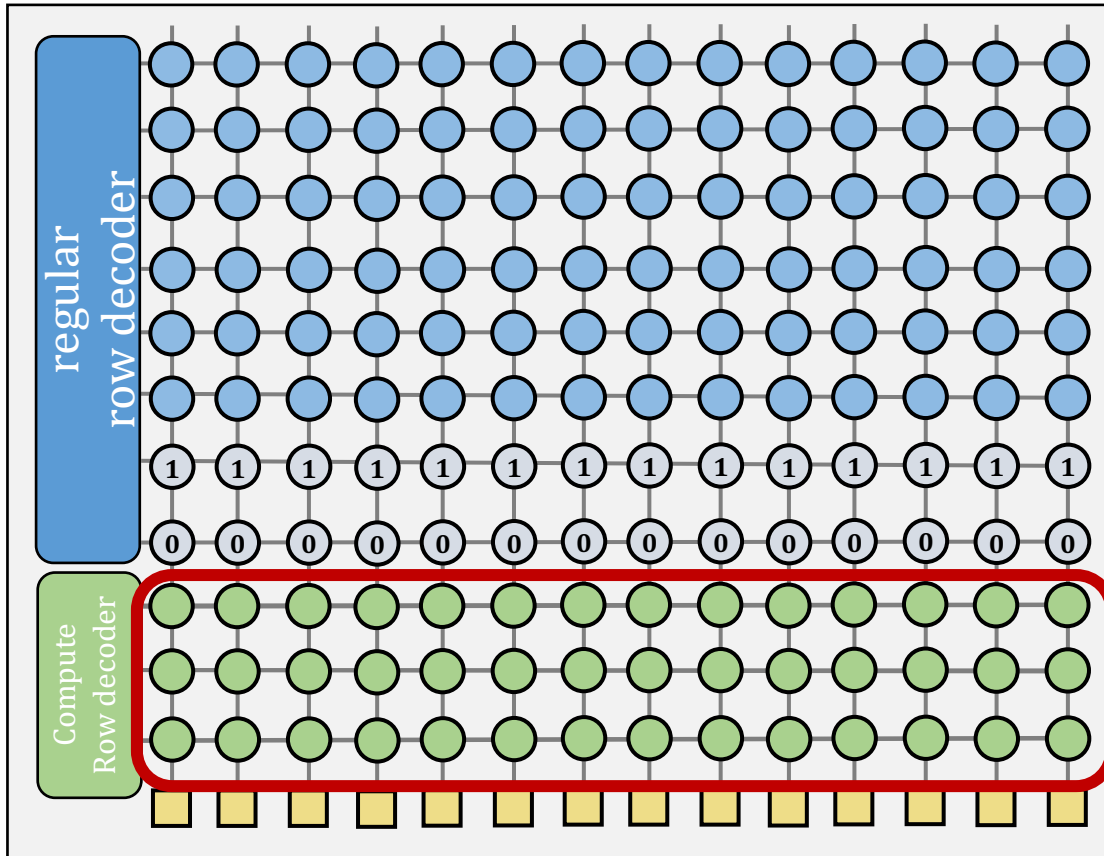
Constraint 1:
Limited number of rows reserved for computation

Compute rows

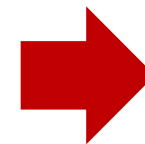
subarray organization

Task 1: Allocating DRAM Rows to Operands

- Allocation algorithm considers **two constraints** specific to processing-using-DRAM



Constraint 2:
Destructive behavior
of triple-row activation

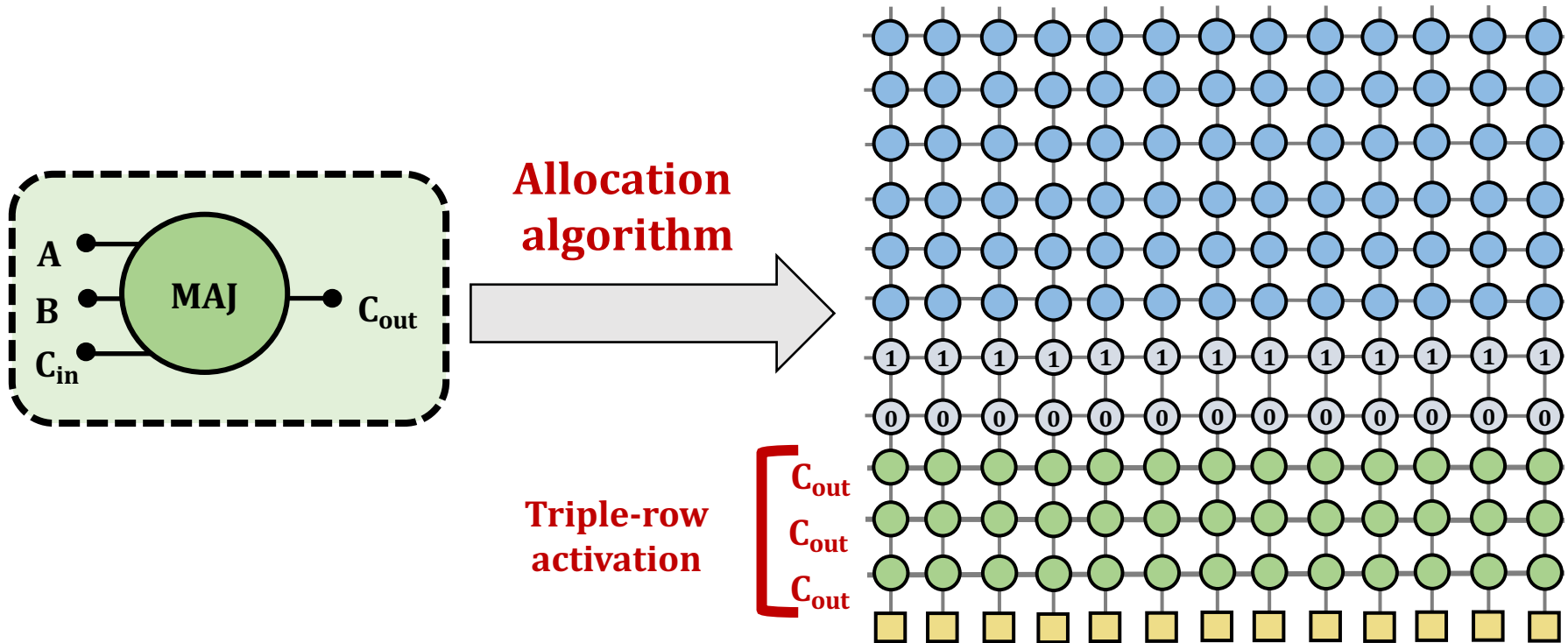


**Overwritten
with MAJ output**

subarray organization

Task 1: Allocating DRAM Rows to Operands

- Allocation algorithm:
 - Assigns as many inputs as the number of **free compute rows**
 - **All three** input rows contain the MAJ output and can be **reused**



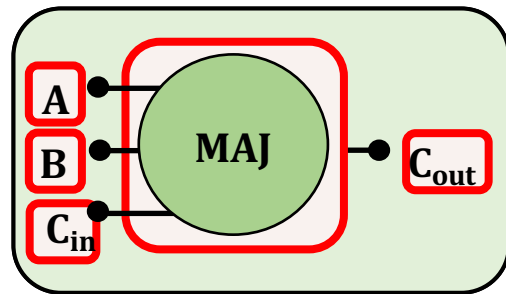
Step 2: μ Program Generation

- **μ Program:** A series of **microarchitectural operations** (e.g., ACT/PRE) that SIMD RAM uses to execute **SIMDRAM operation in DRAM**
- **Goal of Step 2:** To generate the **μ Program** that executes the desired SIMD RAM operation in DRAM

Task 1: Allocate DRAM rows to the operands

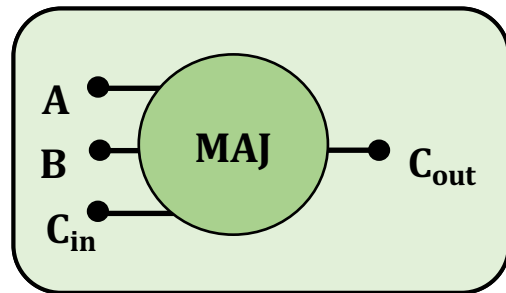
Task 2: Generate μ Program

Task 2: Generate an initial μ Program



1. Generate μ Program

Task 2: Optimize the μ Program



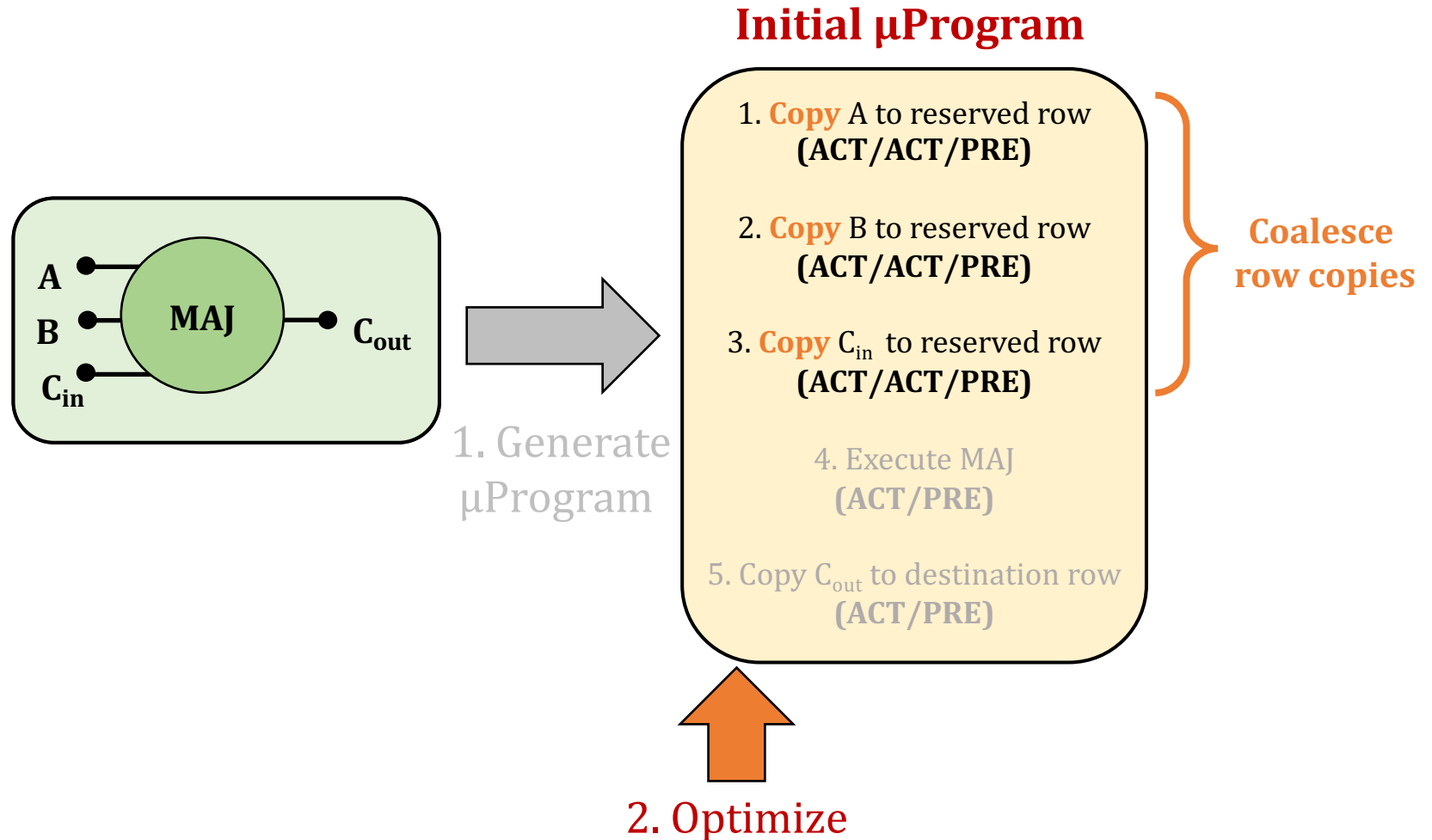
1. Generate μ Program

Initial μ Program

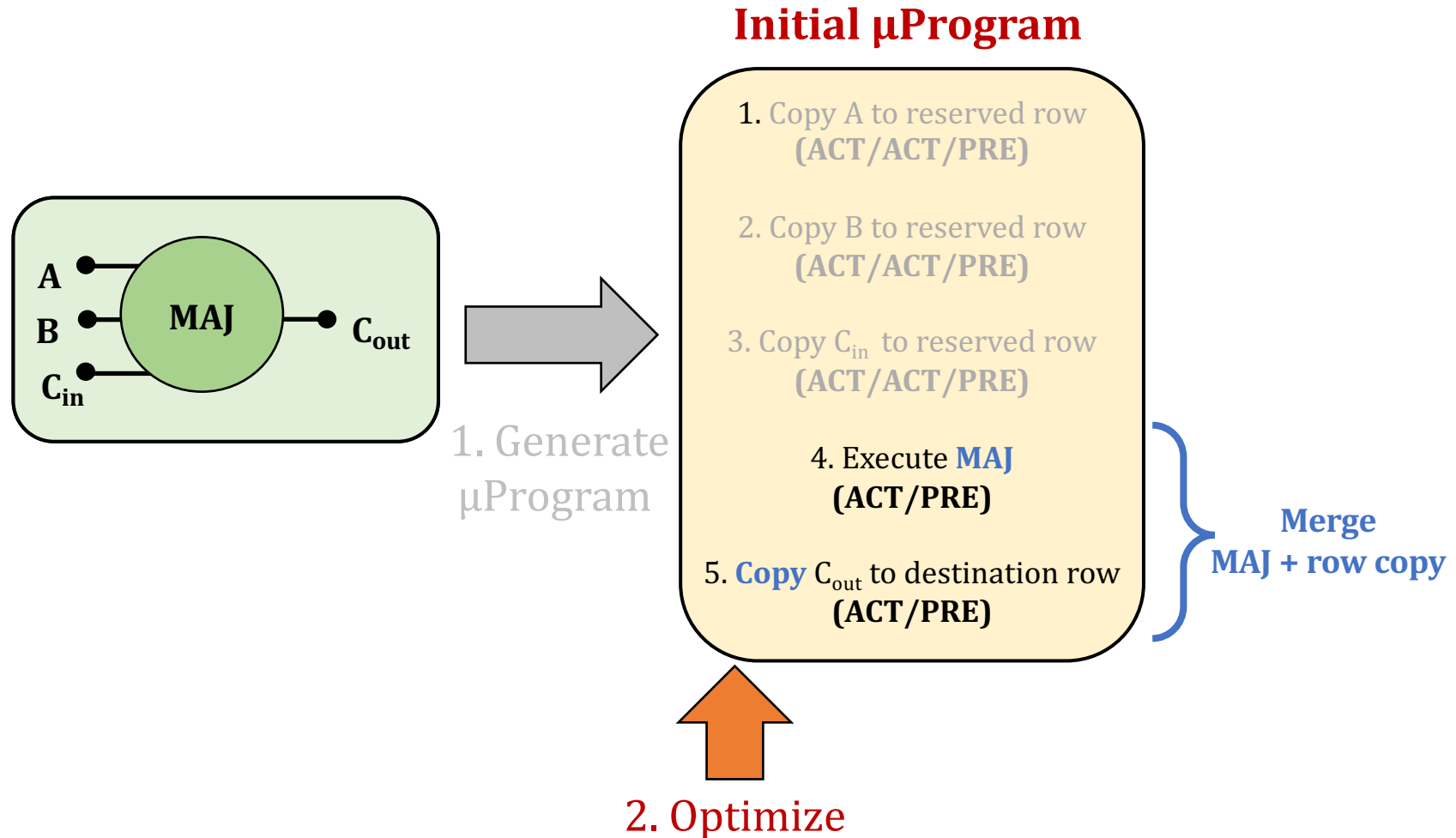
1. Copy A to reserved row
(ACT/ACT/PRE)
2. Copy B to reserved row
(ACT/ACT/PRE)
3. Copy C_{in} to reserved row
(ACT/ACT/PRE)
4. Execute MAJ
(ACT/PRE)
5. Copy C_{out} to destination row
(ACT/PRE)

2. Optimize

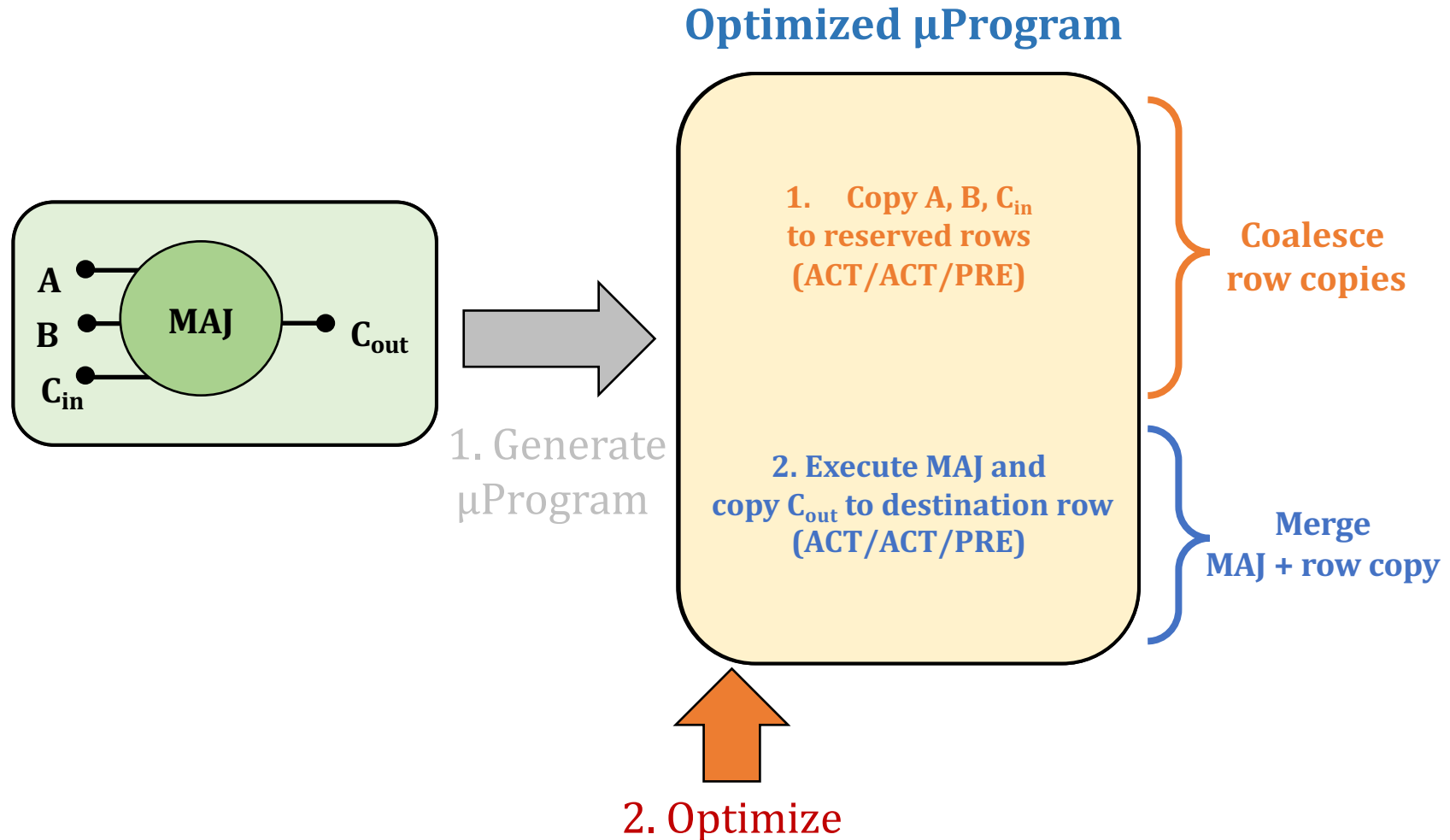
Task 2: Optimize the μ Program



Task 2: Optimize the μ Program

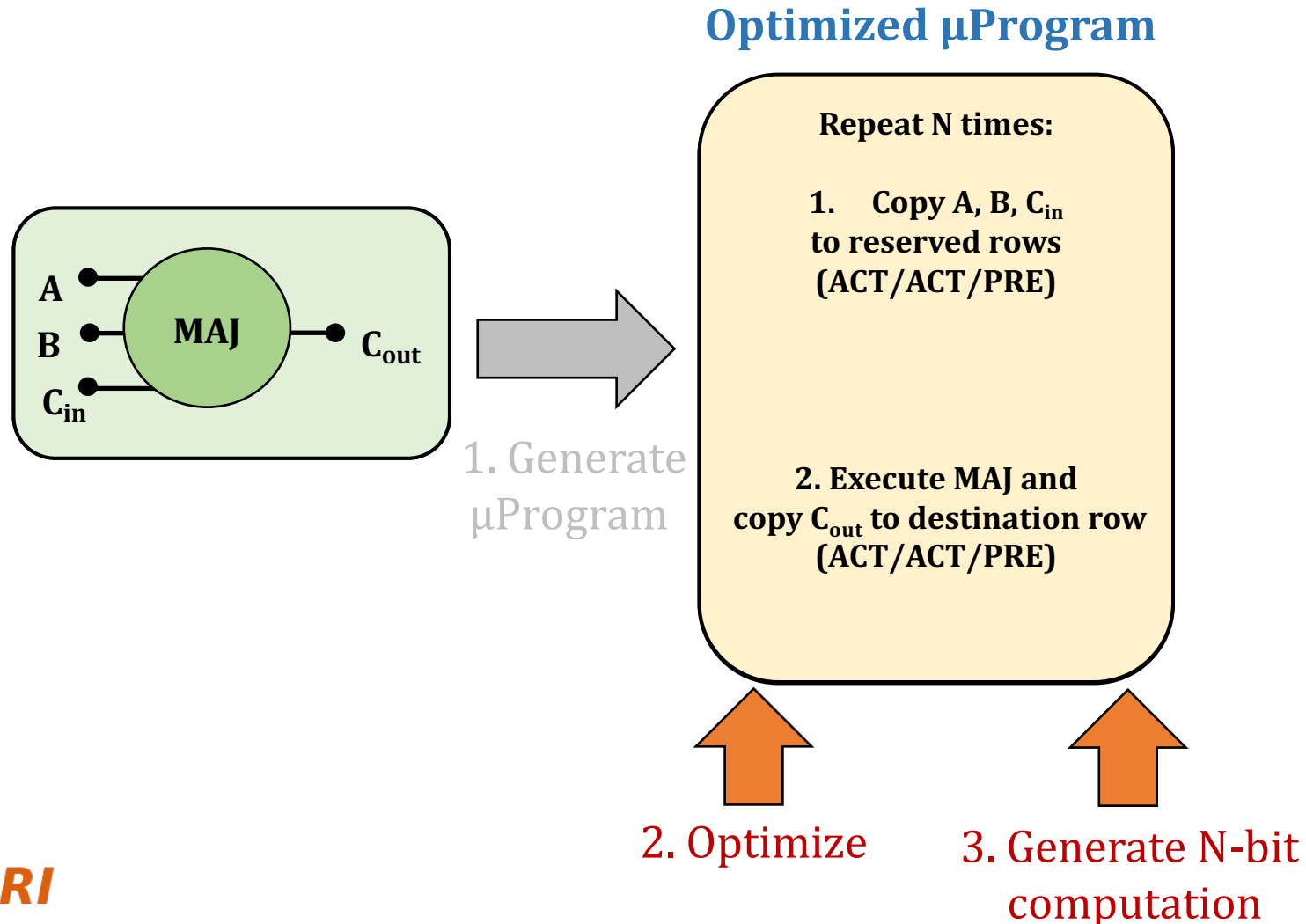


Task 2: Optimize the μ Program



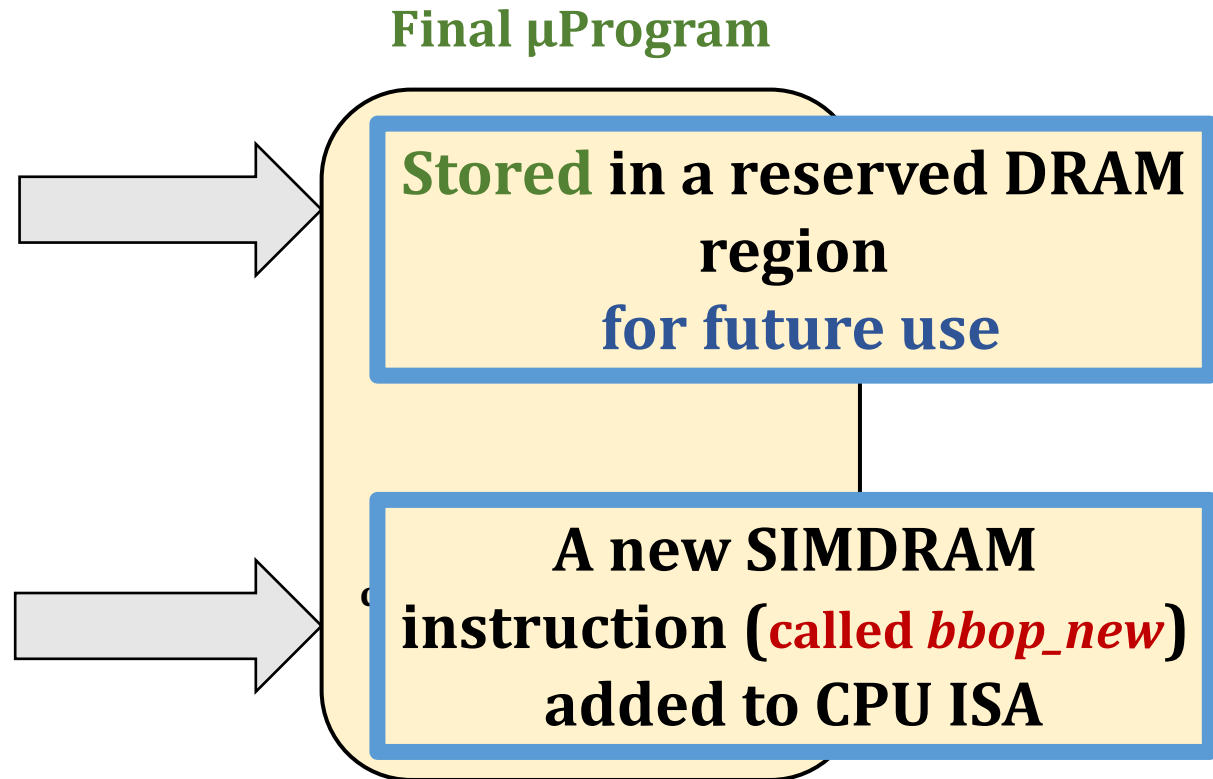
Task 2: Generate N-bit Computation

- **Final μ Program** is optimized and computes the desired operation for operands of N-bit size in a bit-serial fashion

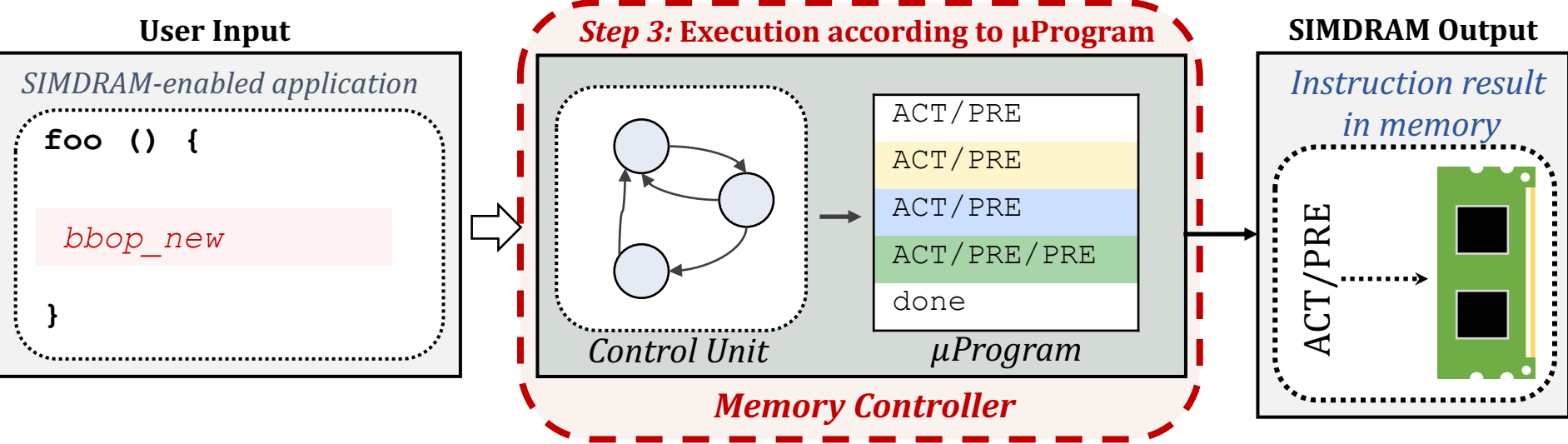
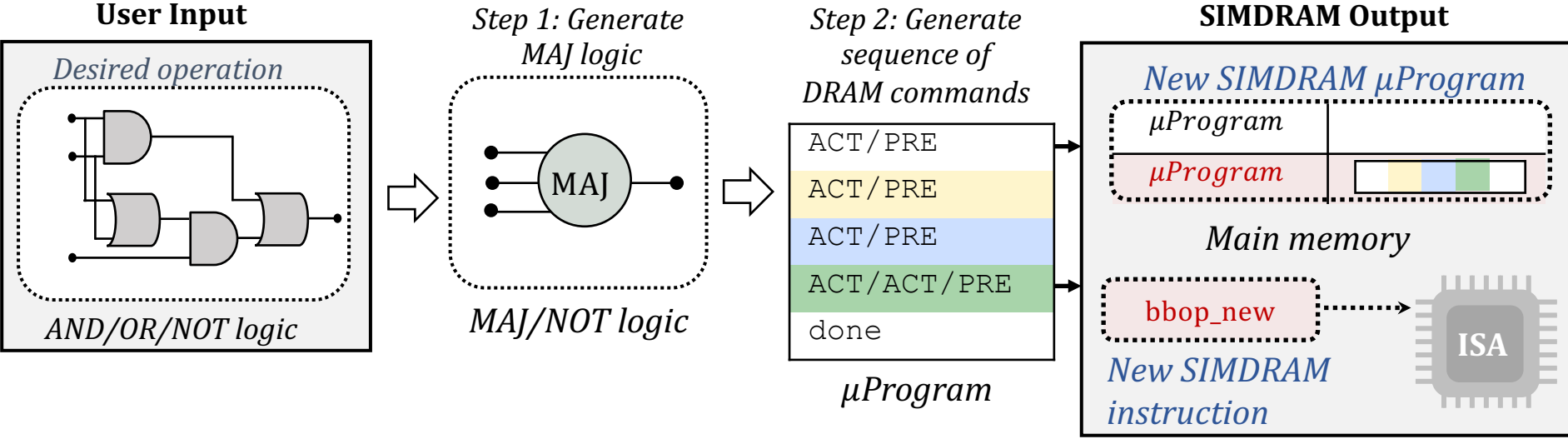


Task 2: Generate μ Program

- **Final μ Program** is optimized and computes the desired operation for operands of N-bit size in a bit-serial fashion

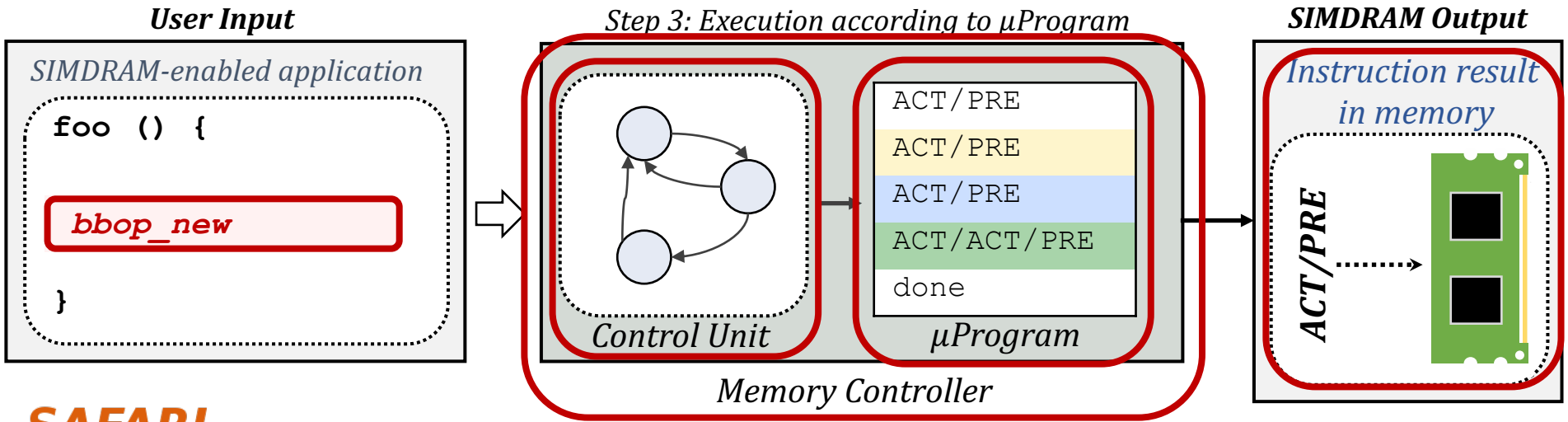


SIMDRAM Framework: Step 3



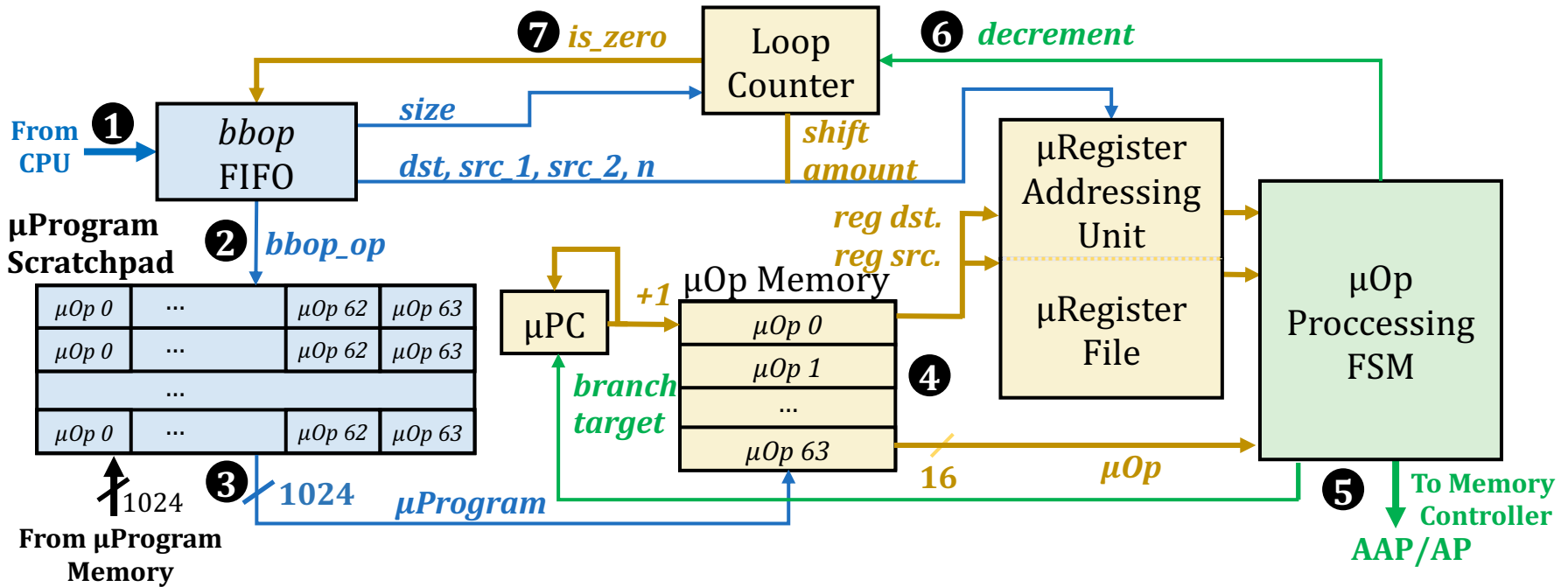
Step 3: μ Program Execution

- **SIMDRAM control unit:** handles the execution of the μ Program at runtime
- Upon receiving a **bbop instruction**, the control unit:
 1. Loads the μ Program corresponding to SIMDRAM operation
 2. Issues the sequence of DRAM commands (ACT/PRE) stored in the μ Program to SIMDRAM subarrays to perform the in-DRAM operation



More in the Paper

- Detailed reference implementation and microarchitecture of the SIMD RAM control unit



Outline

1. Processing-using-DRAM

2. Background

3. SIMD RAM

- Processing-using-DRAM Substrate
- SIMD RAM Framework

4. System Integration

5. Evaluation

6. Conclusion

System Integration

Efficiently transposing data

Programming interface

Handling page faults, address translation,
coherence, and interrupts

Handling limited subarray size

Security implications

Limitations of our framework

System Integration

Efficiently transposing data

Programming interface

Handling page faults, address translation,
coherence, and interrupts

Handling limited subarray size

Security implications

Limitations of our framework

Transposing Data

- SIMD RAM operates on **vertically-laid-out** data
- Other system components expect data to be laid out **horizontally**

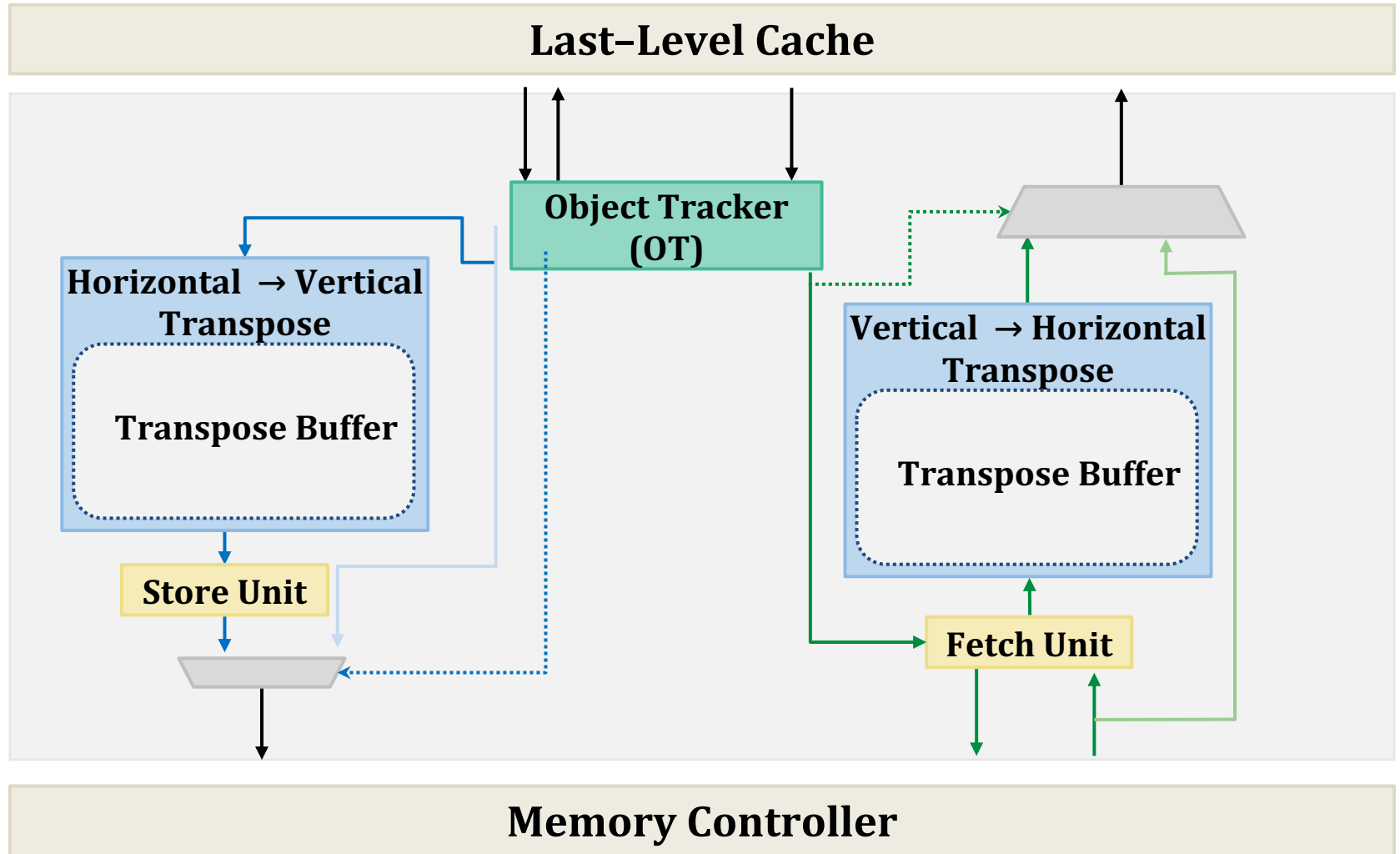


Challenging to share data between SIMD RAM and CPU

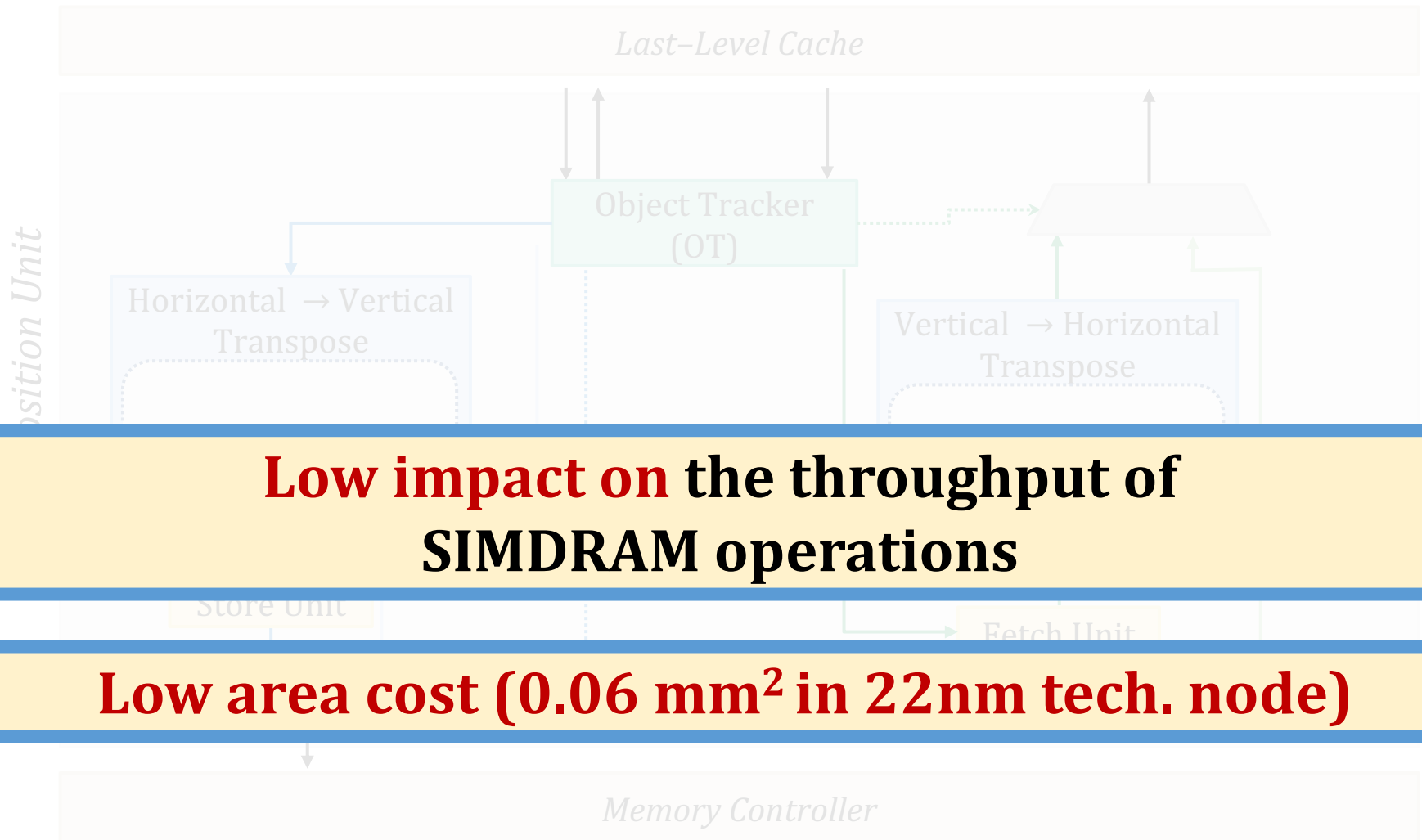
Transposition Unit

Transforms the data layout from **horizontal** to **vertical**, and vice versa

Transposition Unit



Efficiently Transposing Data



Low impact on the throughput of SIMD RAM operations

Low area cost (0.06 mm² in 22nm tech. node)

More in the Paper

SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM

*Nastaran Hajinazar^{1,2} *Geraldo F. Oliveira¹ Sven Gregorio¹ João Dinis Ferreira¹
Nika Mansouri Ghiasi¹ Minesh Patel¹ Mohammed Alser¹ Saugata Ghose³
Juan Gómez-Luna¹ Onur Mutlu¹

¹ETH Zürich

²Simon Fraser University

³University of Illinois at Urbana-Champaign

coherence, and interrupts

Handling limited subarray size

Security implications

Limitations of our framework

Outline

1. Processing-using-DRAM

2. Background

3. SIMD RAM

- Processing-using-DRAM Substrate
- SIMD RAM Framework

4. System Integration

5. Evaluation

6. Conclusion

Methodology: Experimental Setup

- **Simulator:** `gem5`
- **Baselines:**
 - A **multi-core CPU** (Intel Skylake)
 - A **high-end GPU** (NVidia Titan V)
 - **Ambit:** a state-of-the-art in-memory computing mechanism
- **Evaluated SIMD RAM configurations** (all using a DDR4_2400_x64 device):
 - **1-bank:** SIMD RAM exploits 65'536 SIMD lanes (an 8 kB row buffer)
 - **4-banks:** SIMD RAM exploits 262'144 SIMD lanes
 - **16-banks:** SIMD RAM exploits 1'048'576 SIMD lanes

Methodology: Workloads

Evaluated:

- 16 complex in-DRAM operations:

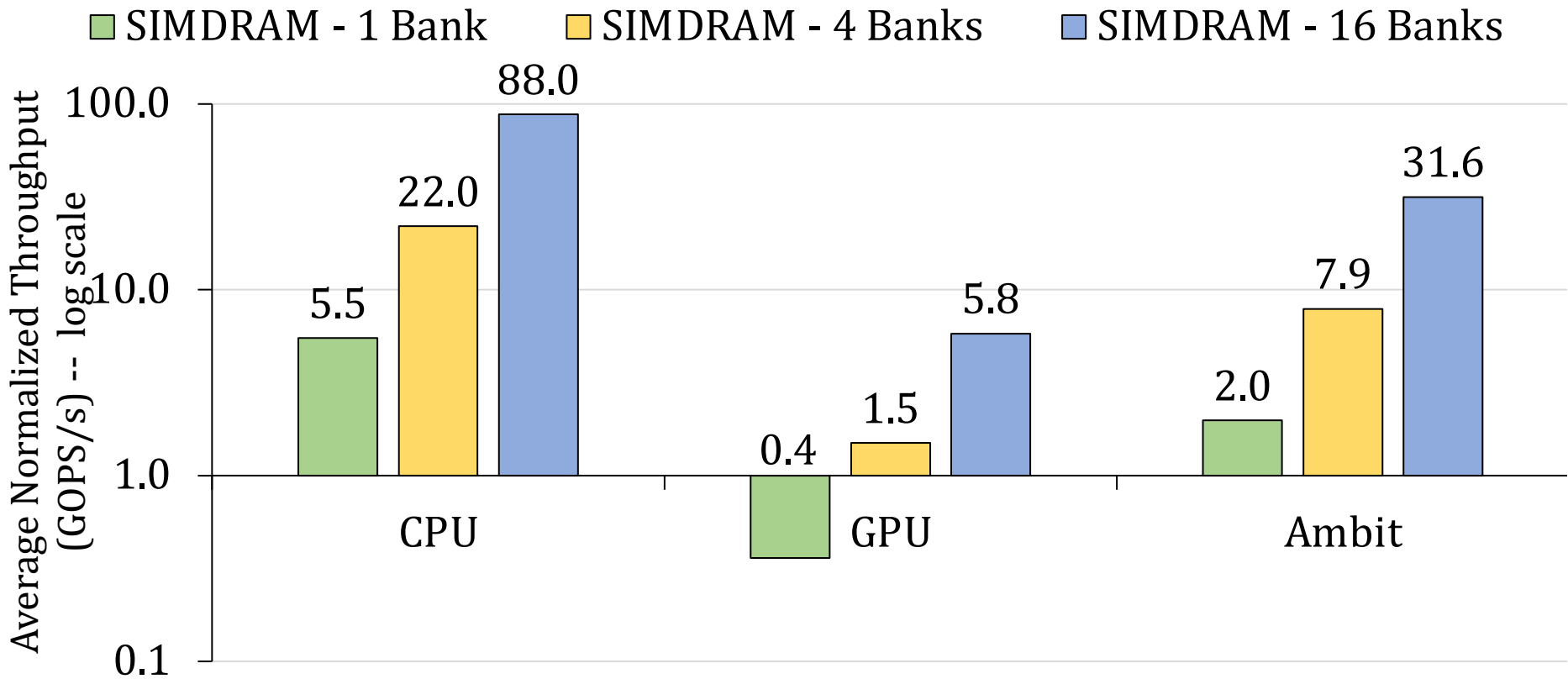
- Absolute
- Addition/Subtraction
- BitCount
- Equality/ Greater/Greater Equal
- Predication
- ReLU
- AND-/OR-/XOR-Reduction
- Division/Multiplication

- 7 real-world applications

- BitWeaving (databases)
- TPH-H (databases)
- kNN (machine learning)
- LeNET (neural networks)
- VGG-13/VGG-16 (neural networks)
- Brightness (graphics)

Throughput Analysis

Average normalized throughput across all 16 SIMD RAM operations

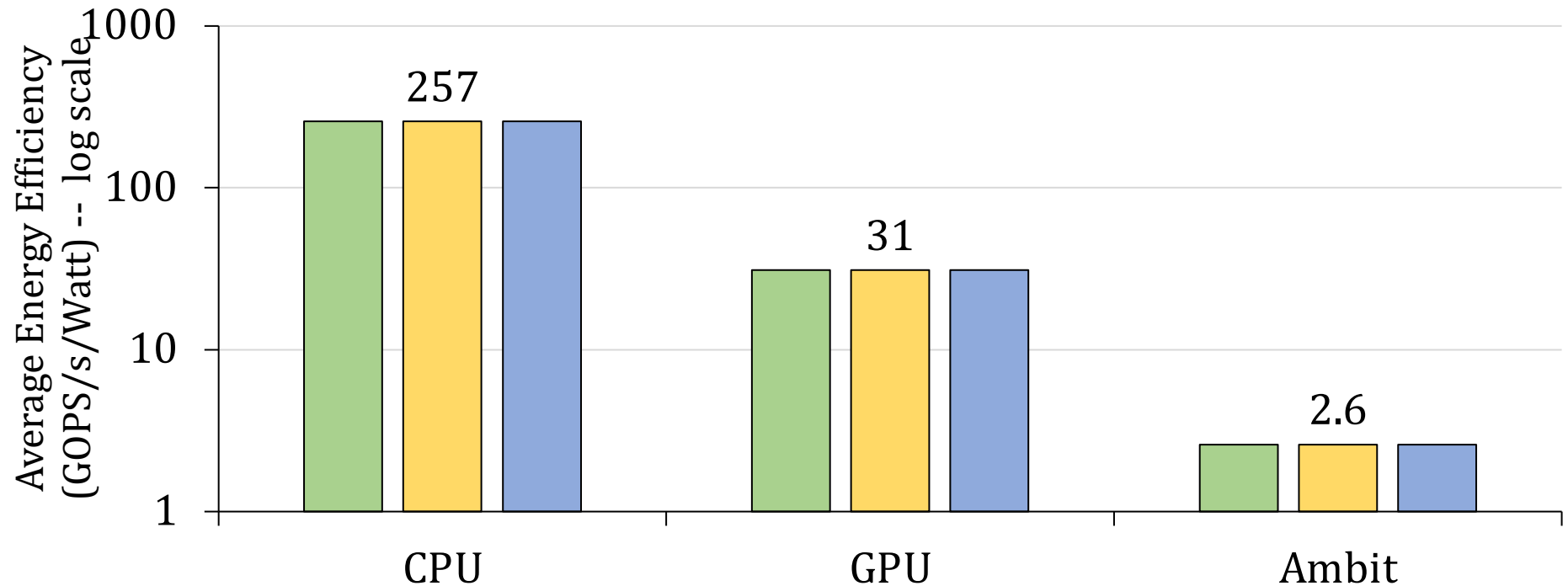


SIMDRAM significantly outperforms
all state-of-the-art baselines for a wide range of operations

Energy Analysis

Average normalized energy efficiency across all 16 SIMD/DRAM operations

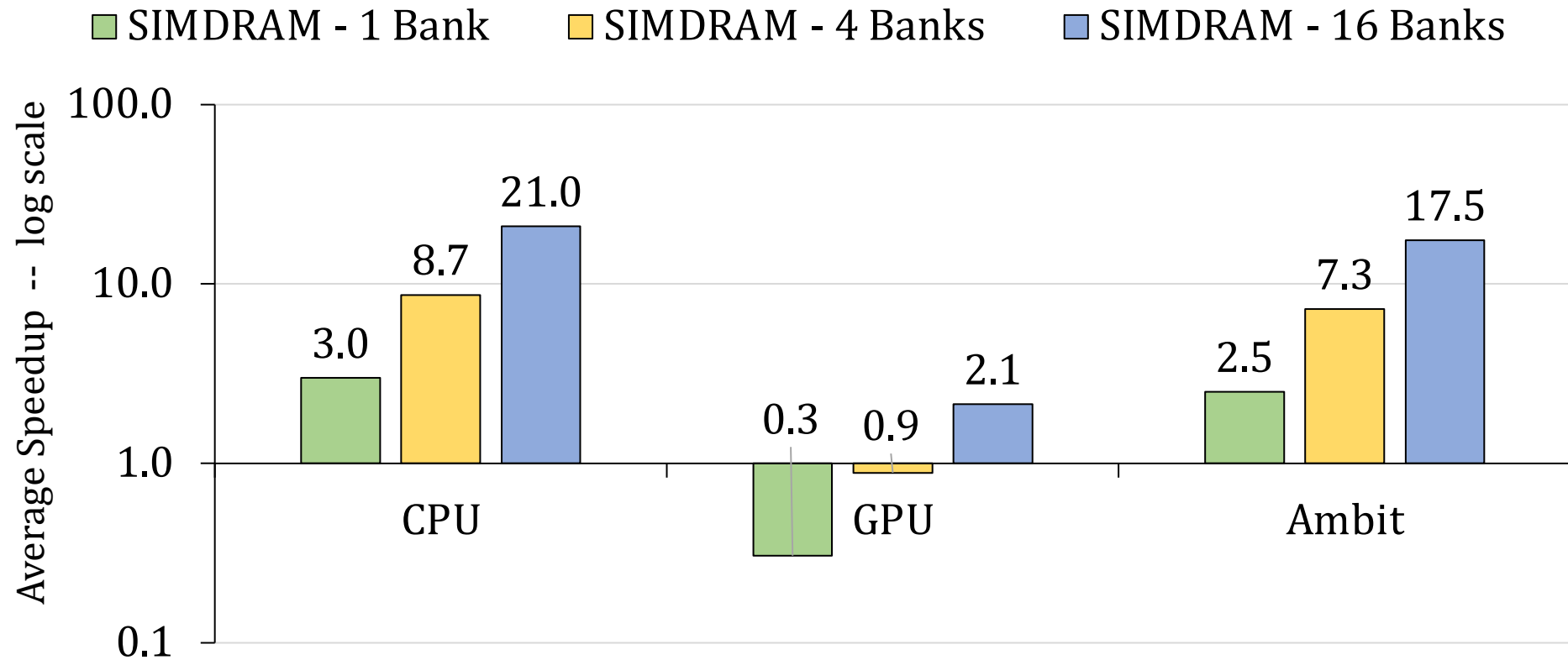
■ SIMD/DRAM - 1 Bank ■ SIMD/DRAM - 4 Banks ■ SIMD/DRAM - 16 Banks



SIMDRAM is more energy-efficient than all state-of-the-art baselines for a wide range of operations

Real-World Applications

Average speedup across 7 real-world applications



SIMDRAM effectively and efficiently accelerates many commonly-used real-world applications

More in the Paper

- **Evaluation:**

- Reliability
- Data transposition overhead
- Area overhead
- Comparison to in-cache computing
- And more ...

Outline

1. Processing-using-DRAM

2. Background

3. SIMDGRAM

- Processing-using-DRAM Substrate
- SIMDGRAM Framework

4. System Integration

5. Evaluation

6. Conclusion

Conclusion

- **SIMDRAM**: An end-to-end processing-using-DRAM framework that provides the programming interface, the ISA, and the hardware support for:
 1. **Efficiently** computing complex operations
 2. Providing the ability to implement **arbitrary** operations as required
 3. Using a **massively-parallel** in-DRAM SIMD substrate
- **Key Results**: SIMDRAM provides:
 - **88x** and **5.8x** the **throughput** and **257x** and **31x** the **energy efficiency** of a baseline CPU and a high-end GPU, respectively, for 16 in-DRAM operations
 - **21x** and **2.1x** the **performance** of the CPU and GPU for seven real-world applications
- **Conclusion**: SIMDRAM is a promising PuM framework
 - Can **ease the adoption** of processing-using-DRAM architectures
 - Improve the **performance** and **efficiency** of processing-using-DRAM architectures

SIMDRAM: A Framework for Bit-Serial SIMD Processing using DRAM

Nastaran Hajinazar*

Geraldo F. Oliveira*

Sven Gregorio

Joao Ferreira

Nika Mansouri Ghiasi

Minesh Patel

Mohammed Alser

Saugata Ghose

Juan Gómez-Luna

Onur Mutlu

SAFARI

ETH zürich

