

A New MRAM-based Process In-Memory Accelerator for Efficient Neural Network Training with Floating Point Precision



Hongjie Wang, Yang Zhao*, Chaojian Li, Yue Wang, and Yingyan Lin*

Department of Electrical and Computer Engineering, Rice University

2020 IEEE International Symposium on Circuits and Systems
Virtual, October 10-21, 2020



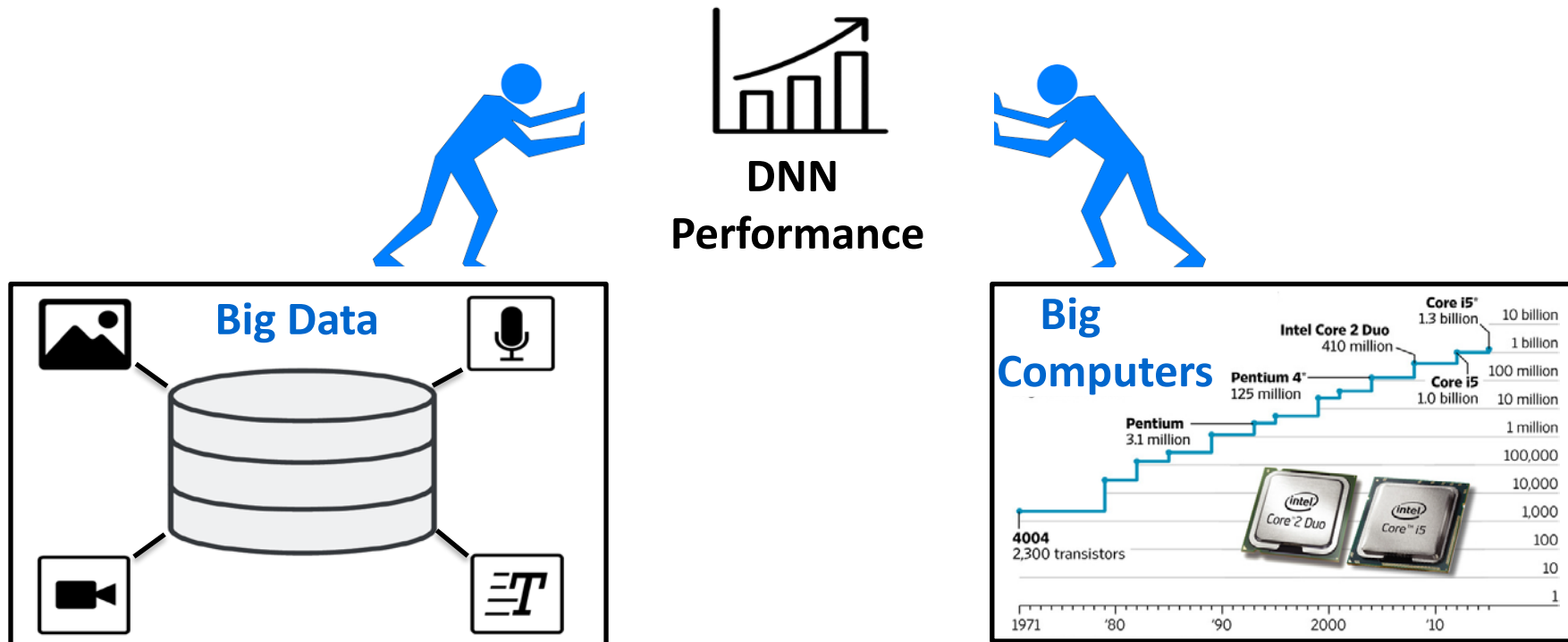
*the first two authors contribute equally

Outline

- **Background and Motivation**
- ***The Proposed Accelerator: Overview***
- ***The Proposed Accelerator: Computing Methodology***
- ***Evaluation of the Proposed Accelerator***
- **Conclusion**

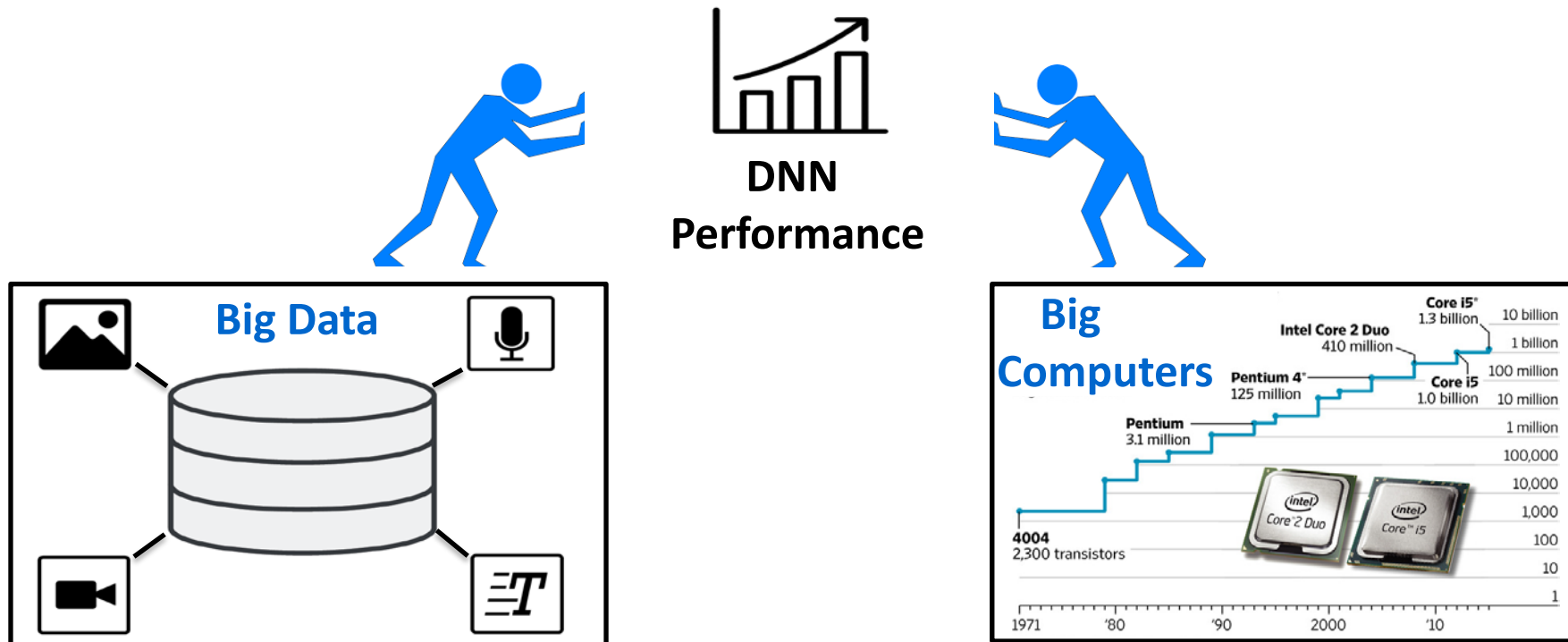
Growing Demand for Efficient DNN Training

- Powerful DNNs **require** training with **big data + big computers**



Growing Demand for Efficient DNN Training

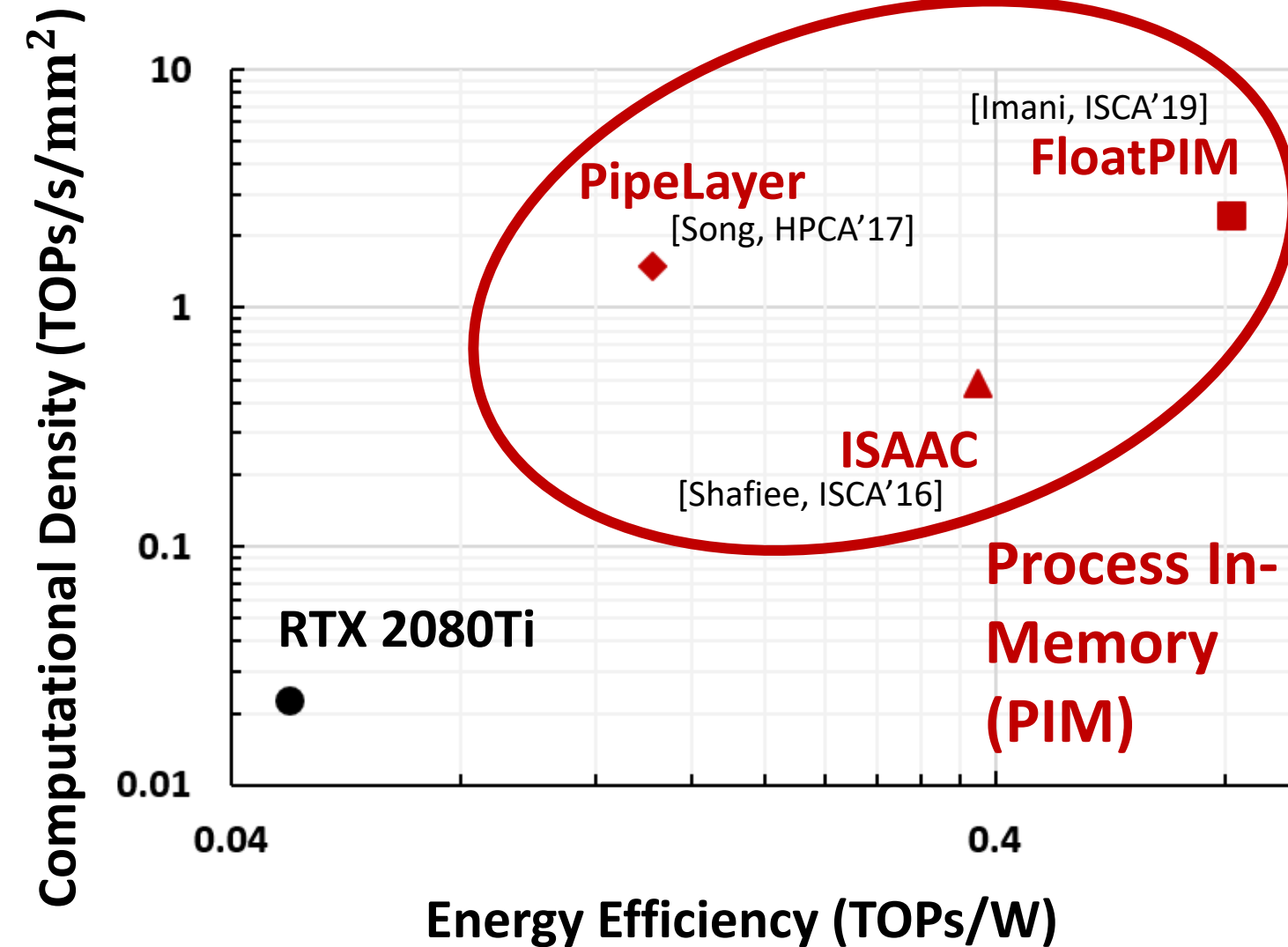
- Powerful DNNs **require** training with **big data + big computers**



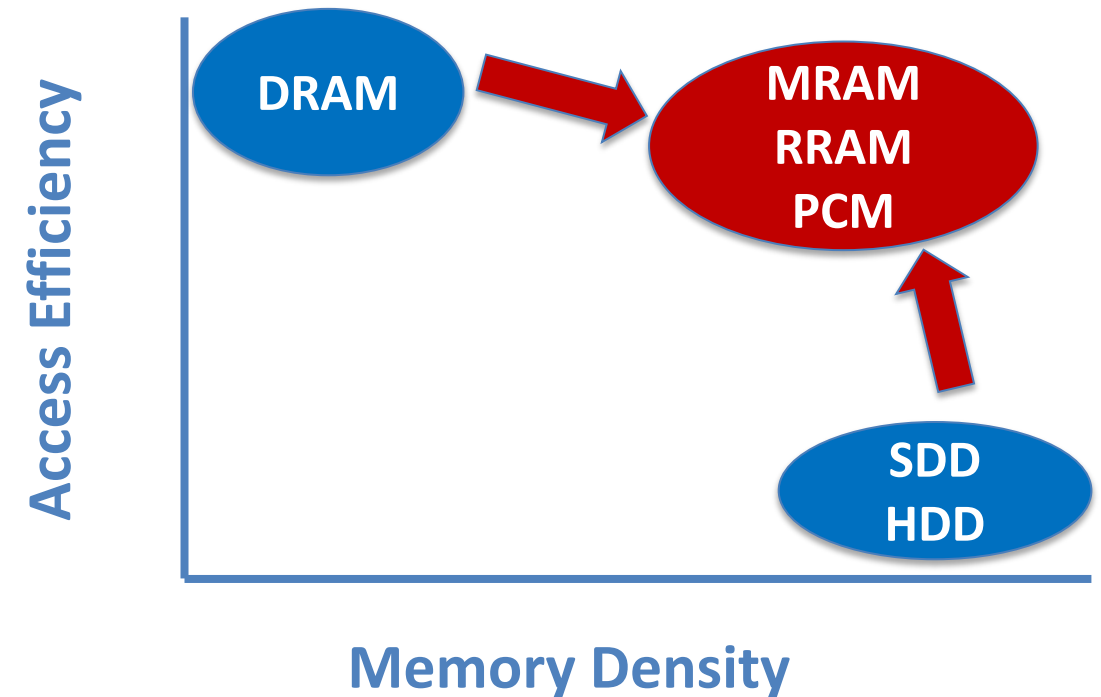
- DNN training challenges:**
 - Prohibitive computational and time cost**
(e.g., 10^{18} FLOPs for training ResNet50 ~ 14 days using a NVIDIA M40 GPU)
 - Increasing environmental concern**
(e.g., carbon emission of training a DNN \approx one car's lifetime emission)

Process In-Memory for DNN Acceleration

- Powerful DNNs require training with **big data** + big computers

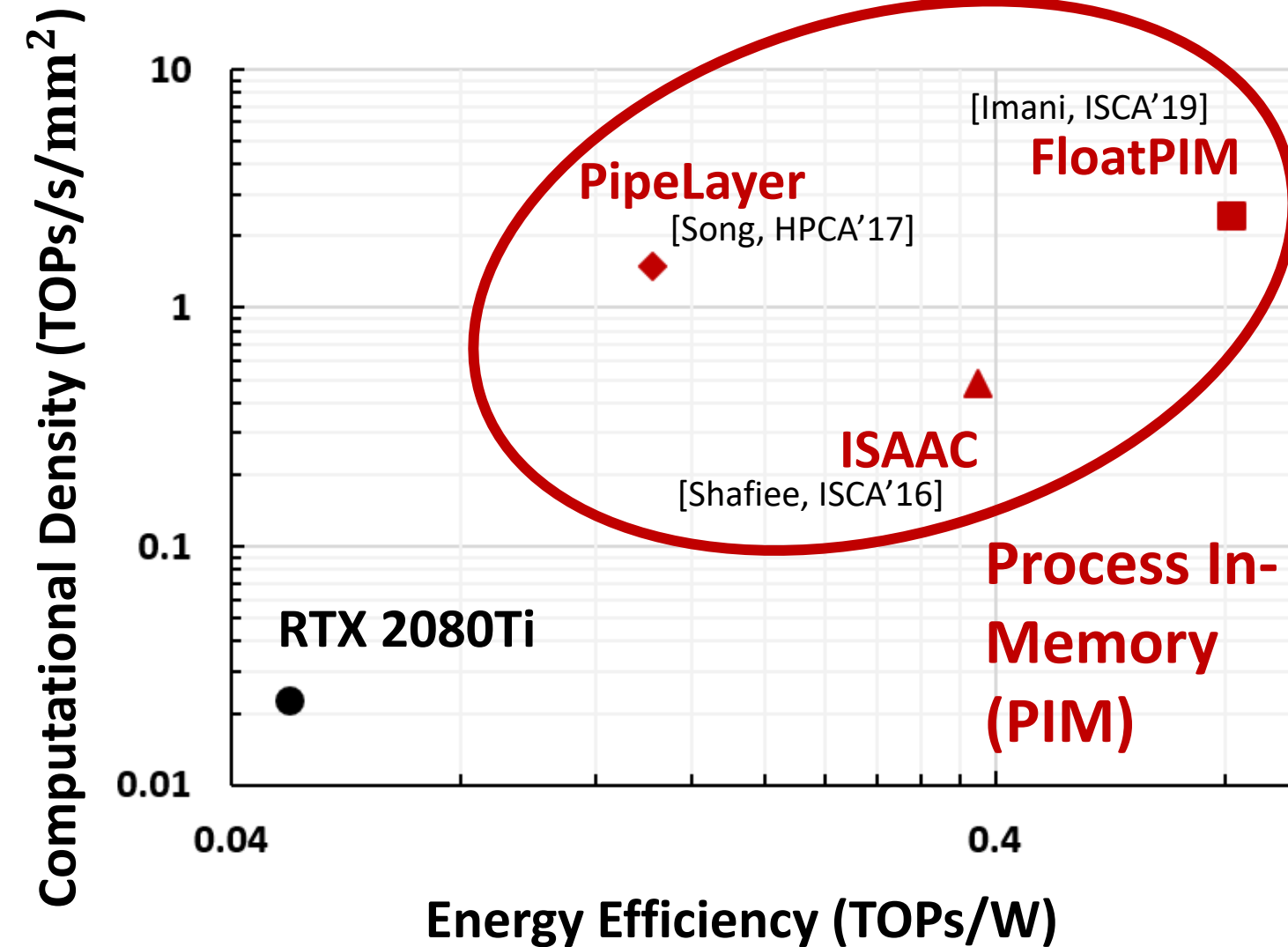


- High density non-volatile memory for huge number of parameters

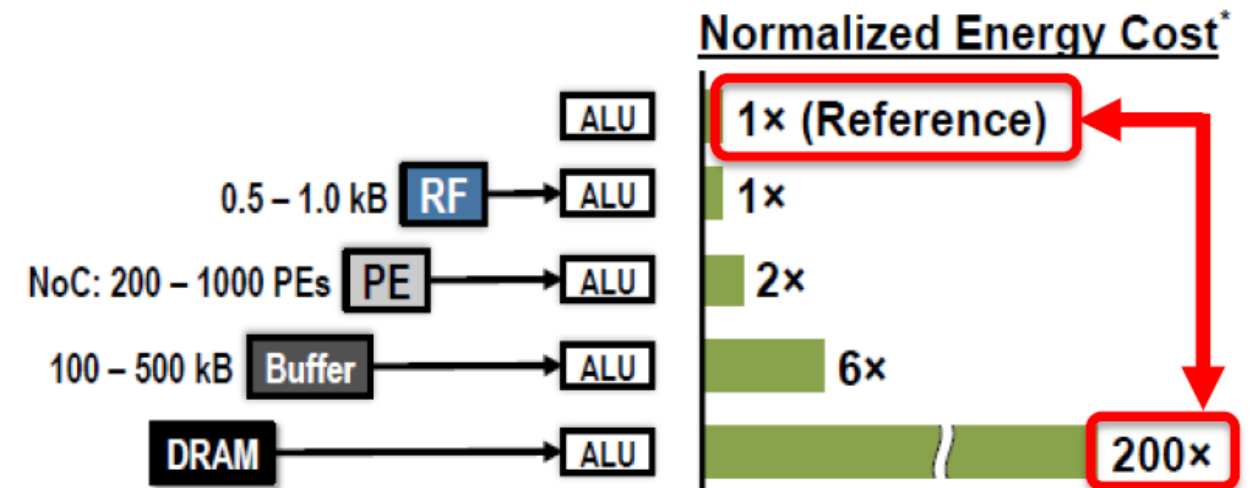


Process In-Memory for DNN Acceleration

- Powerful DNNs require training with big data + **big computers**



- High density non-volatile memory for huge number of parameters
- Decrease costly data movements in Von Neumann Architecture



* measured from a commercial 65nm process
[Chen, ISCA'16]

Existing PIM Accelerators for DNN Training

- **Existing PIM Accelerators for Deep Neural Network**

- ☹️ Support **only inference** with low precision [Shafiee, ISCA'16] [Chi, ISCA'16][Patil, ISCAS'19]
- ☹️ **Costly writing** for widely used ReRAM based accelerator [Imani, ISCA'19] [Song, HPCA'17]
- ☹️ **High computational complexity** caused by limited available Boolean functions [Imani, ISCA'19]

- **Our Proposed PIM Accelerator**

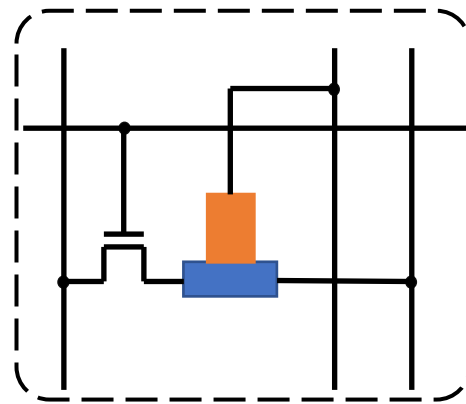
- 😊 Supports **both inference and training** with floating point precision
- 😊 **Low writing current** thanks to SOT-MRAM based implementation
- 😊 **More efficient computing** caused by plentiful available Boolean functions

Outline

- Background and Motivation
- *The Proposed Accelerator: Overview*
- *The Proposed Accelerator: Computing Methodology*
- Evaluation of *the Proposed Accelerator*
- Conclusion

The Proposed Accelerator: Overview

Read and Write: Cell Design

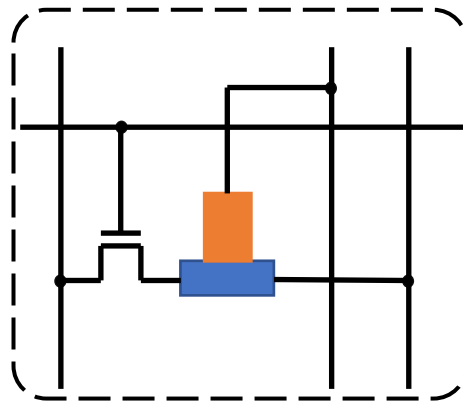


1T-1R
structure

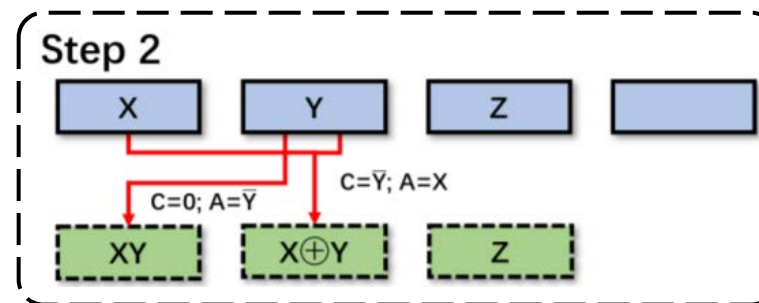
The Proposed Accelerator: Overview

1-Bit Full Addition: Cell Level Data Transfer

Read and Write: Cell Design



1T-1R
structure



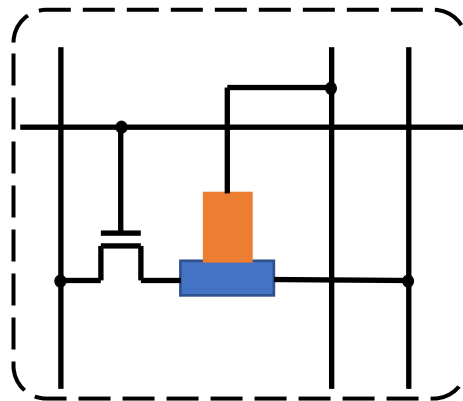
4 steps to perform
1-bit full addition

The Proposed Accelerator: Overview

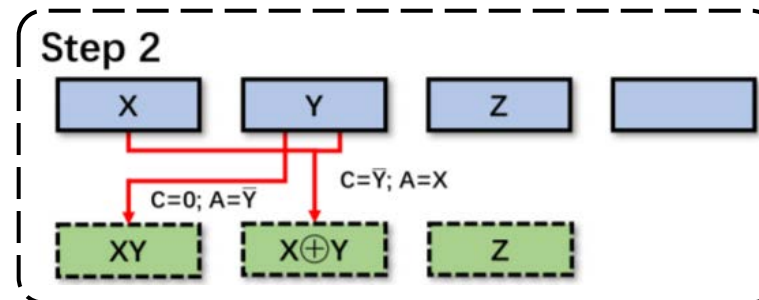
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

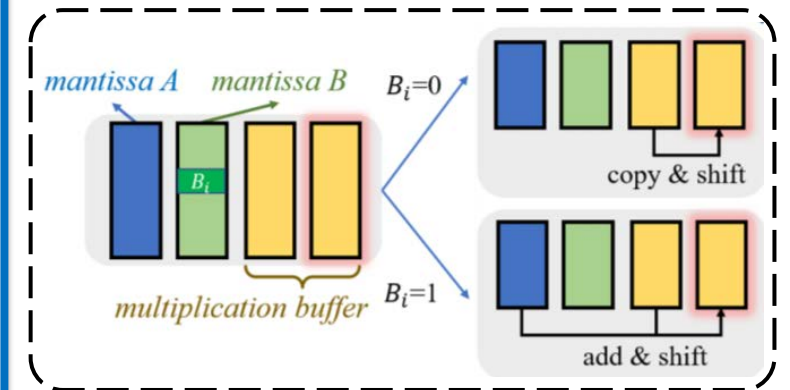
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Addition and multiplication
design

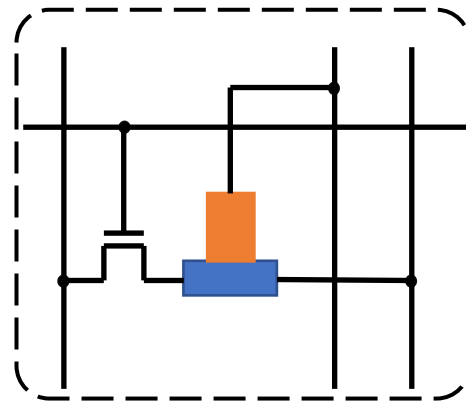
The Proposed Accelerator: Overview

Training and Inference: *the Proposed Accelerator*

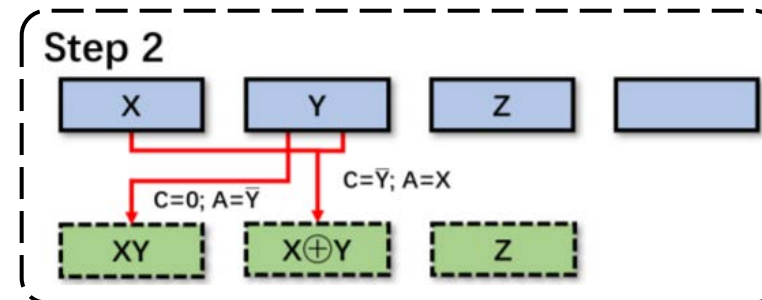
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

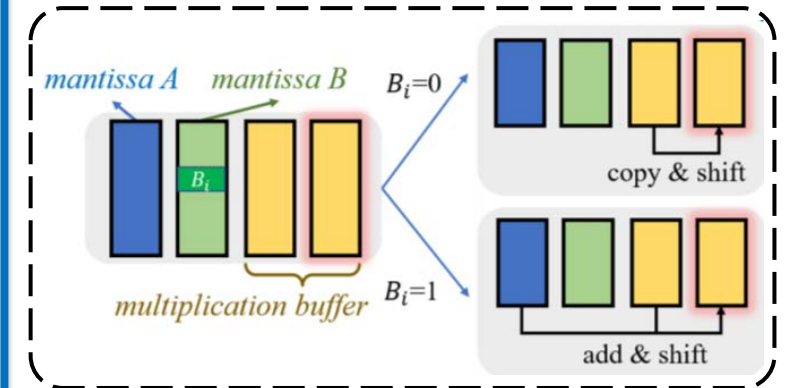
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition

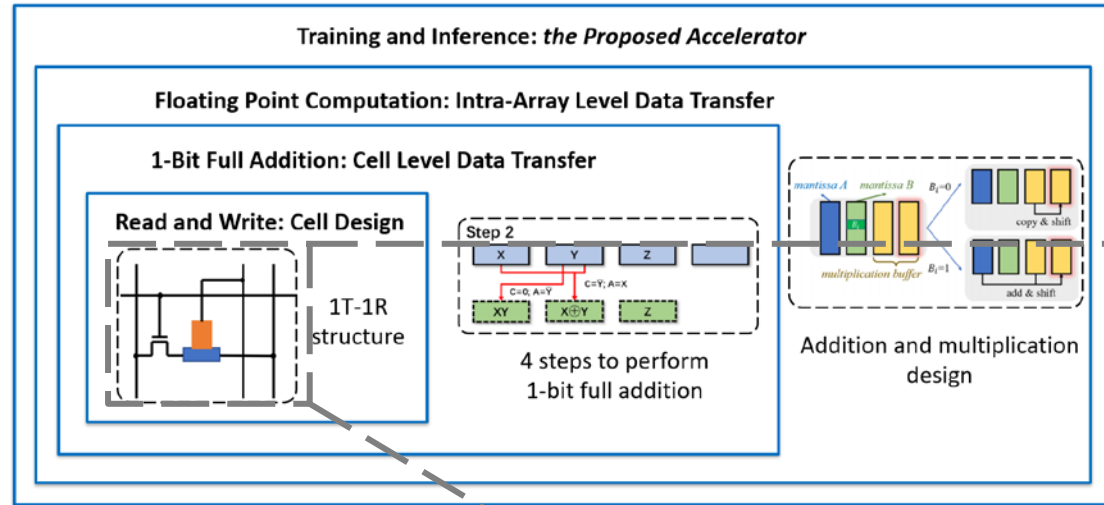


Addition and multiplication
design

Outline


- Background and Motivation
- *The Proposed Accelerator: Overview*
- *The Proposed Accelerator: Computing Methodology*
- Evaluation of *the Proposed Accelerator*
- Conclusion

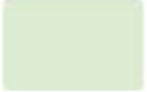
Contribution 1: SOT-MRAM Cell Design

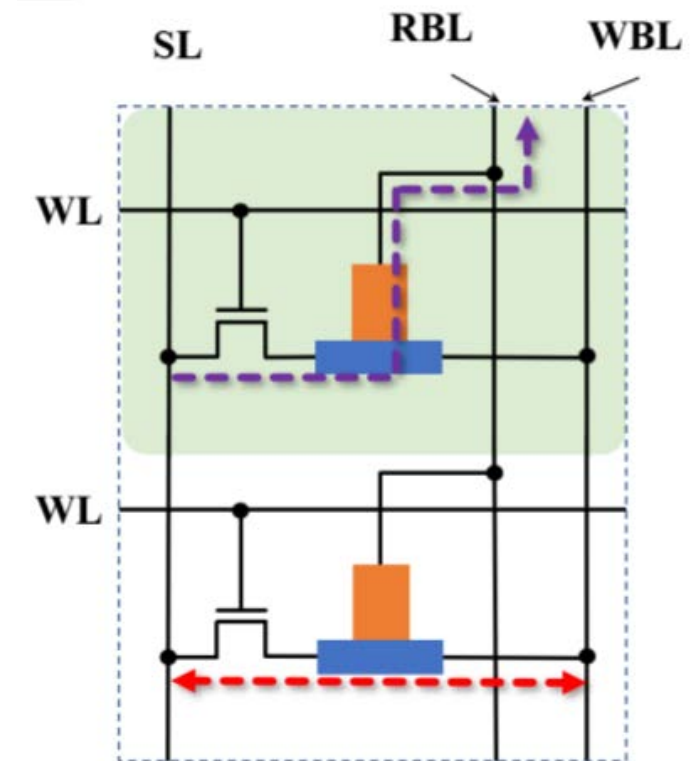


Read current 
Write current 

MTJ 

Heavy metal 

One memory cell 



1T-1R (*proposed*)

- High flexibility
- Fair density

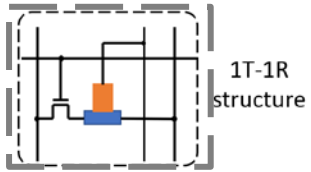
Contribution 1: SOT-MRAM Cell Design

Training and Inference: *the Proposed Accelerator*

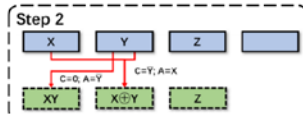
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

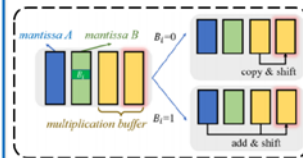
Read and Write: Cell Design



1T-1R structure



4 steps to perform 1-bit full addition



Addition and multiplication design

Read current 
Write current 

MTJ



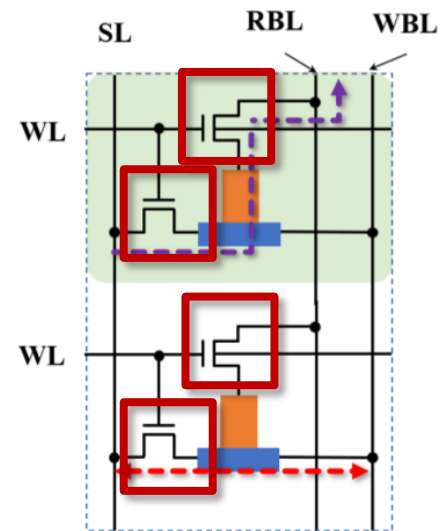
Heavy metal



One memory cell



transistor in cell



2T-1R

- High flexibility
- Low density

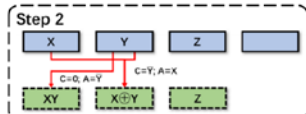
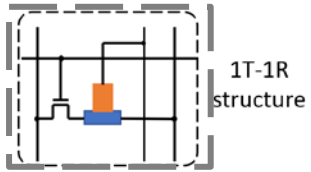
Contribution 1: SOT-MRAM Cell Design

Training and Inference: *the Proposed Accelerator*

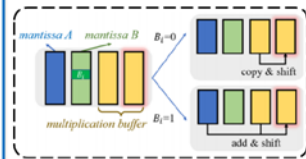
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

Read and Write: Cell Design





4 steps to perform 1-bit full addition



Read current 
Write current 

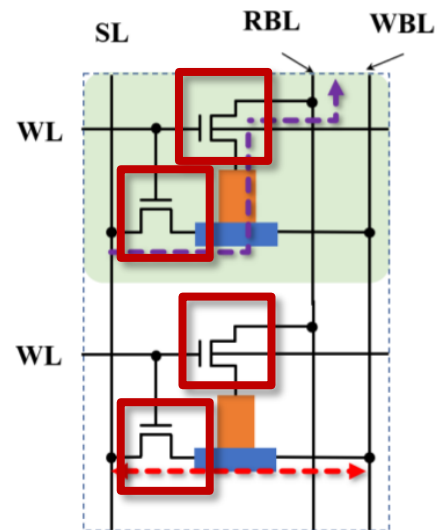
MTJ 

Heavy metal 

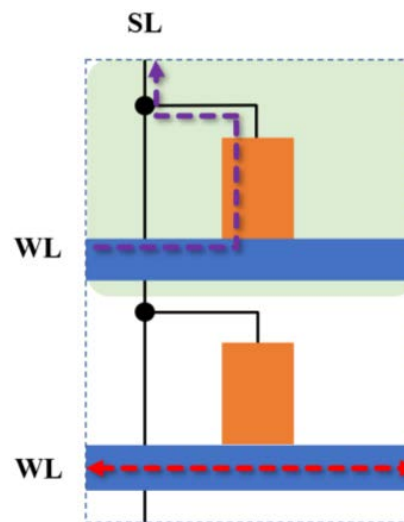
One memory cell 



transistor in cell

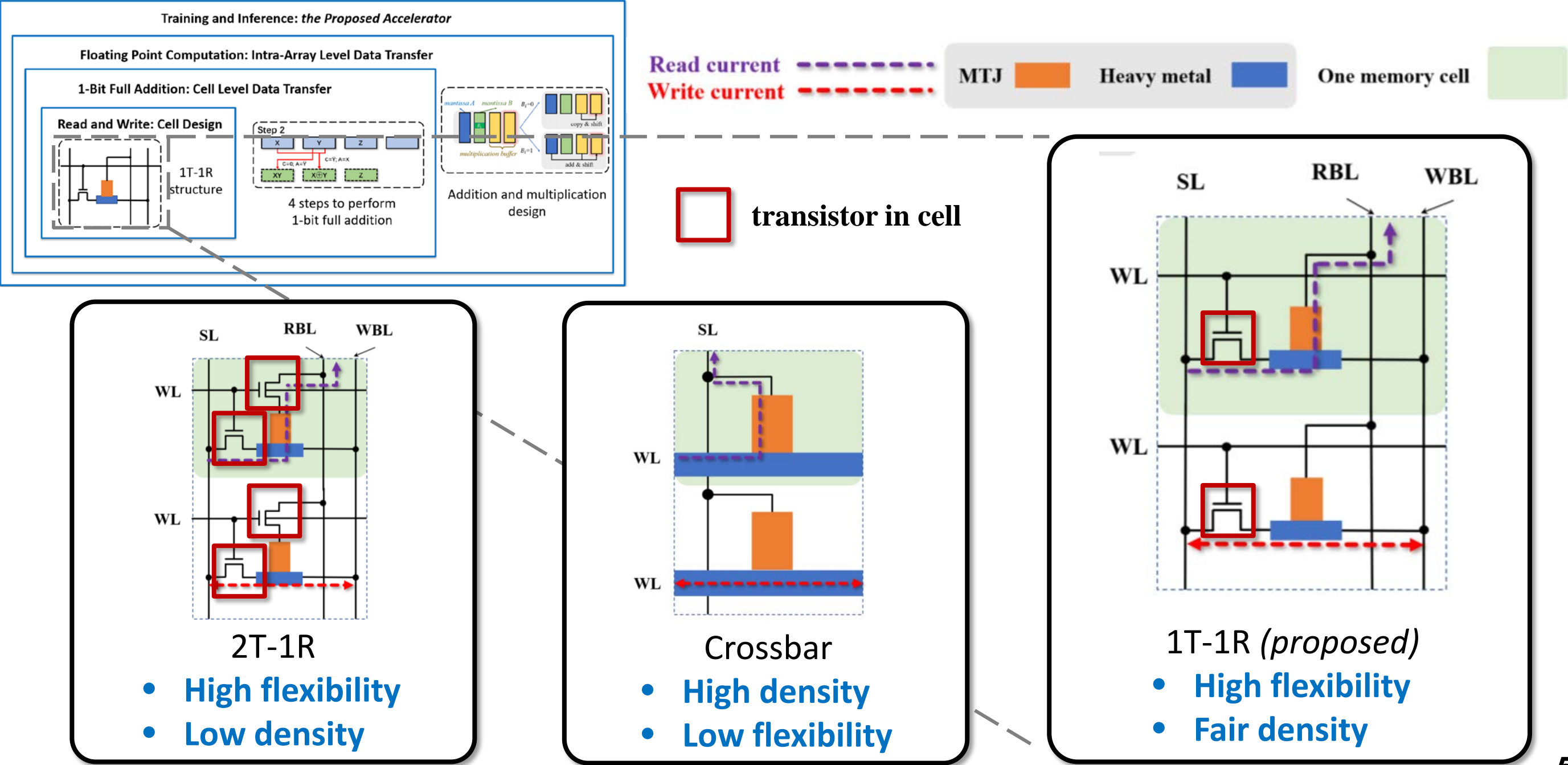


- High flexibility
- Low density

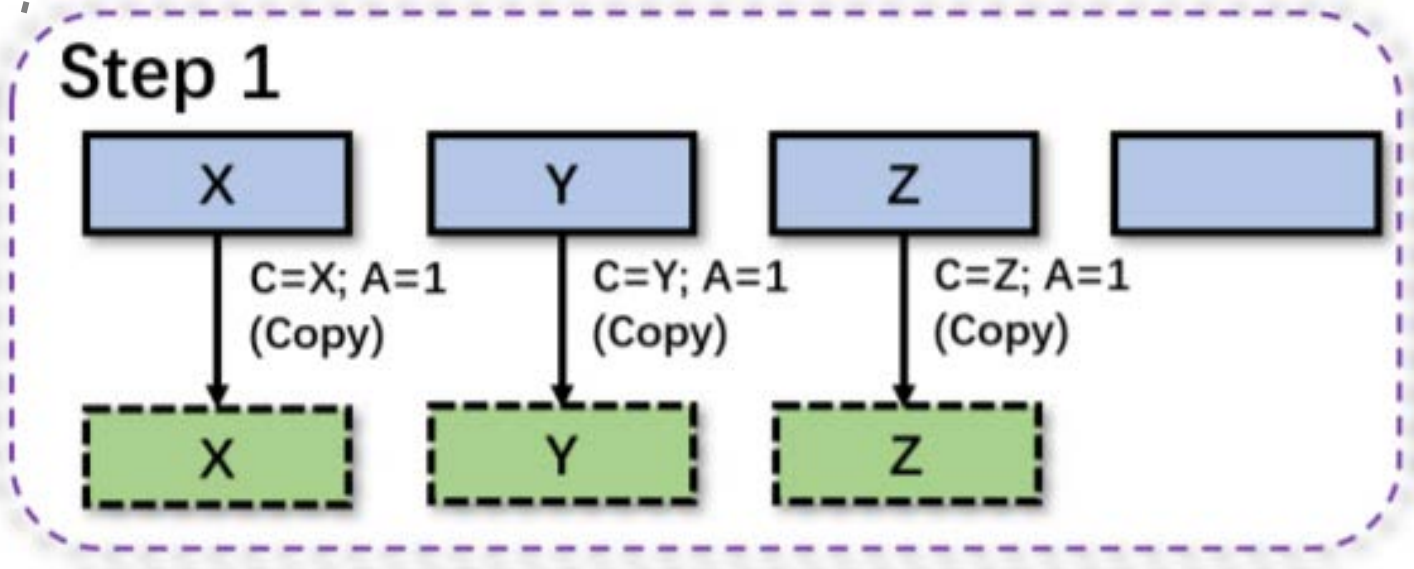
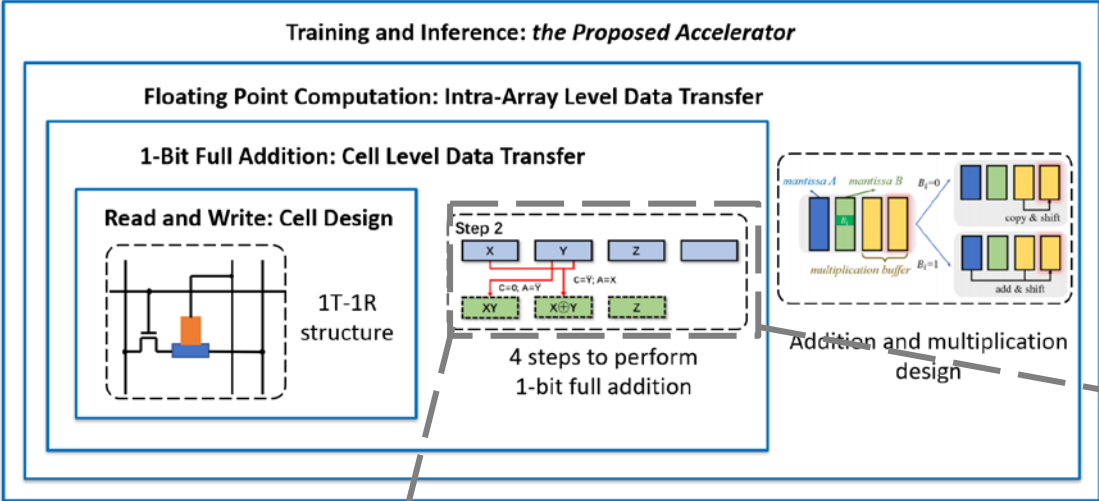


- High density
- Low flexibility

Contribution 1: SOT-MRAM Cell Design



Contribution 2: Full Adder Based on the 1T-1R Cell



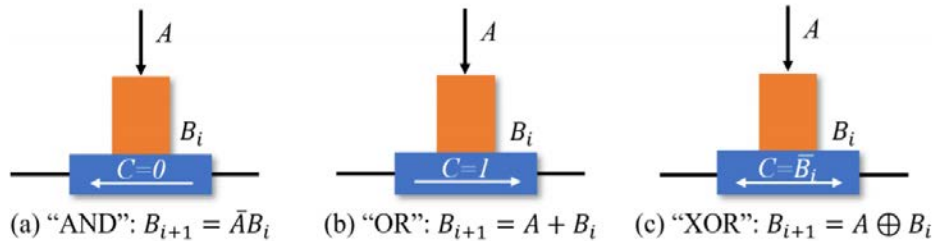
MRAM

MRAM cache

→ Data movement

→ Computation

Definition of A , B and C



logic functions in the SOT-MRAM cell

[Zhang, T-ED'17]

X, Y : operands
 Z : input carry
 Z' : output carry
 $S(R)$: sum

$$S(R) = X \oplus Y \oplus Z$$
$$Z' = XY + Z(X \oplus Y)$$

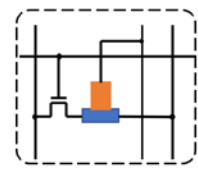
Contribution 2: Full Adder Based on the 1T-1R Cell

Training and Inference: *the Proposed Accelerator*

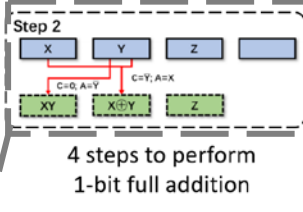
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

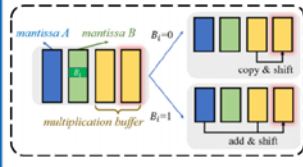
Read and Write: Cell Design



1T-1R structure

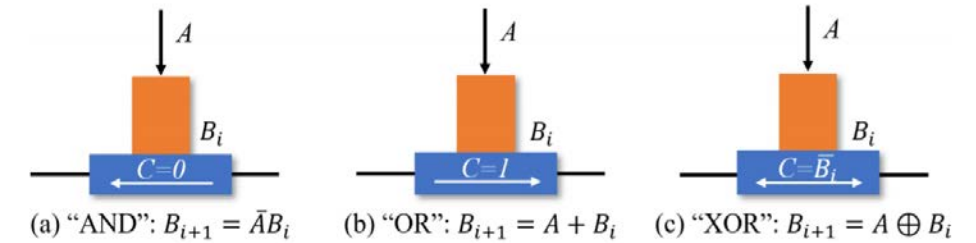


4 steps to perform 1-bit full addition



Addition and multiplication design

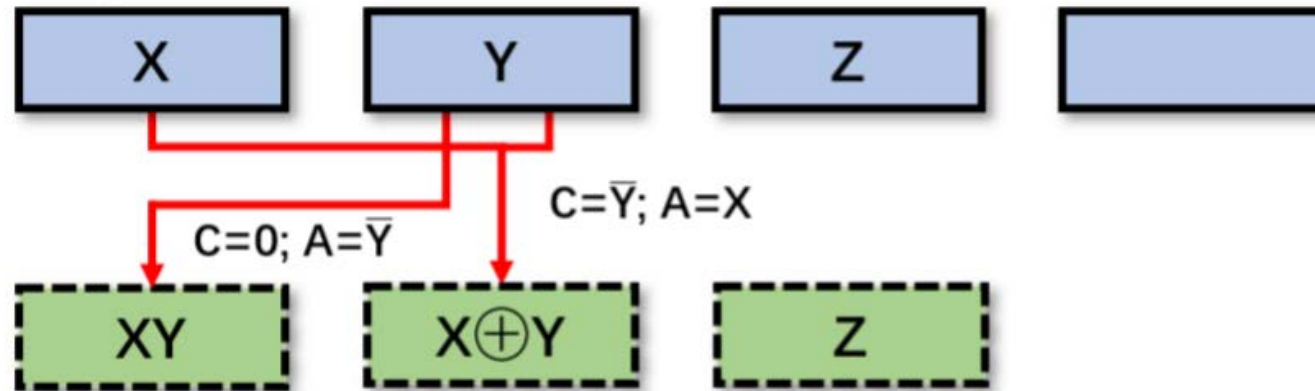
Definition of A , B and C



logic functions in the SOT-MRAM cell

[Zhang, T-ED'17]

Step 2



X, Y : operands
 Z : input carry
 Z' : output carry
 $S(R)$: sum

$$S(R) = X \oplus Y \oplus Z$$

$$Z' = XY + Z(X \oplus Y)$$



MRAM

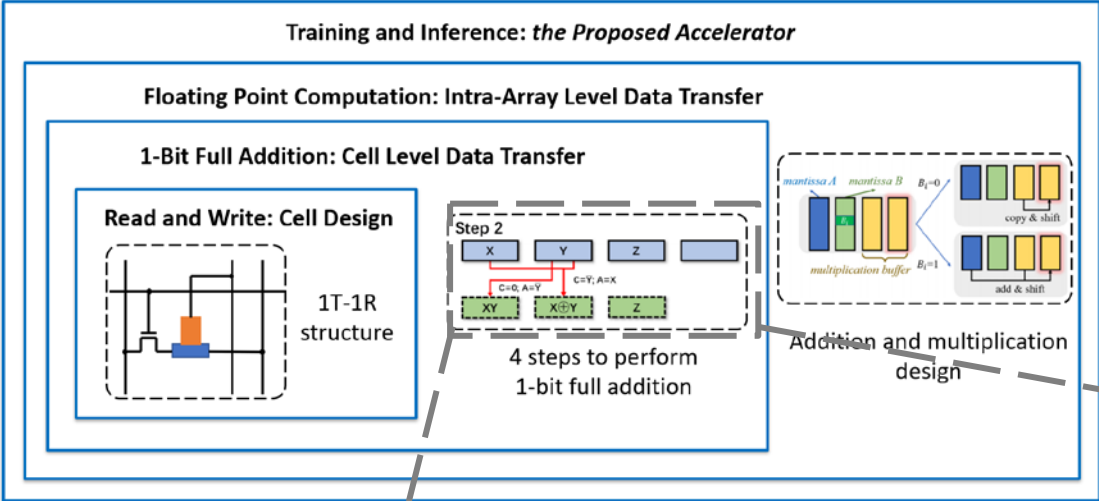


MRAM cache

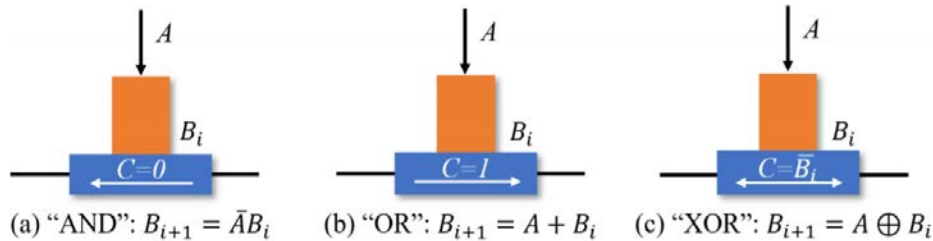
→ Data movement

→ Computation

Contribution 2: Full Adder Based on the 1T-1R Cell

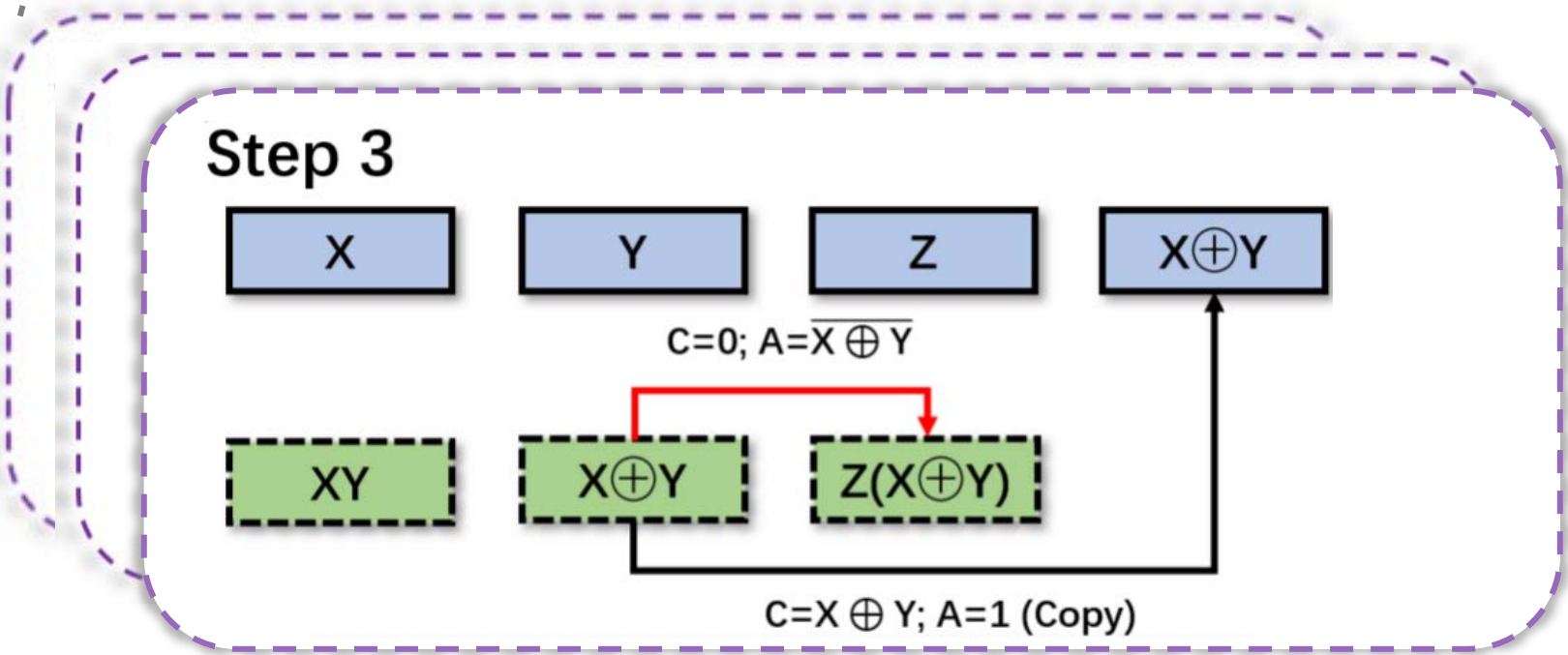


Definition of A , B and C



logic functions in the SOT-MRAM cell

[Zhang, T-ED'17]



X, Y : operands
 Z : input carry
 Z' : output carry
 $S(R)$: sum

$$S(R) = X \oplus Y \oplus Z$$
$$Z' = XY + Z(X \oplus Y)$$



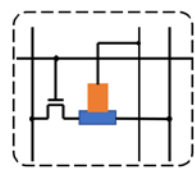
Contribution 2: Full Adder Based on the 1T-1R Cell

Training and Inference: *the Proposed Accelerator*

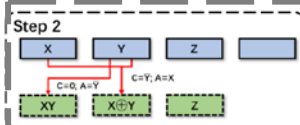
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

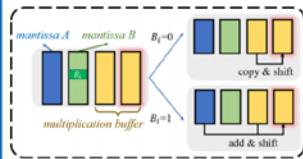
Read and Write: Cell Design



1T-1R structure

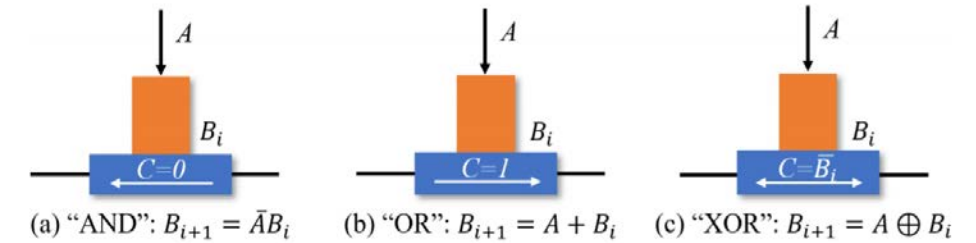


4 steps to perform 1-bit full addition



Addition and multiplication design

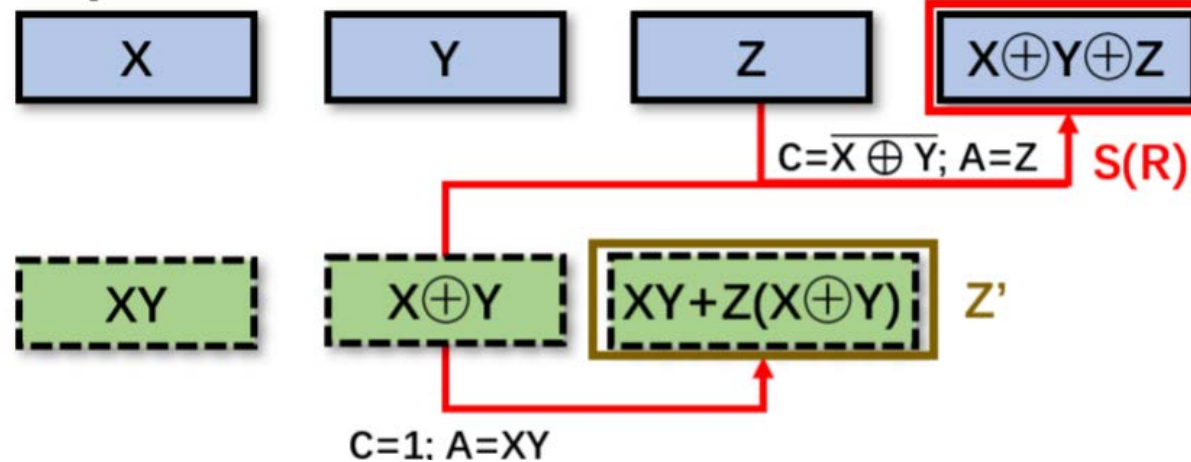
Definition of A , B and C



logic functions in the SOT-MRAM cell

[Zhang, T-ED'17]

Step 4



X, Y : operands
 Z : input carry
 Z' : output carry
 $S(R)$: sum

$$S(R) = X \oplus Y \oplus Z$$

$$Z' = XY + Z(X \oplus Y)$$



MRAM



MRAM cache

→ Data movement

→ Computation

Complexity of the Proposed Full Adder Over SOTA One

Compare Complexity in Terms of Required Clock Cycles and Memory Cells

PIM Implementation	Read and Write Clock Cycles	Extra Memory Cells
The Proposed Method	4	3
FloatPIM [Imani, ISCA'19]	13	10

Reduced by:

69.2%↓

70.0%↓

Available Boolean Function vs. Computational Complexity

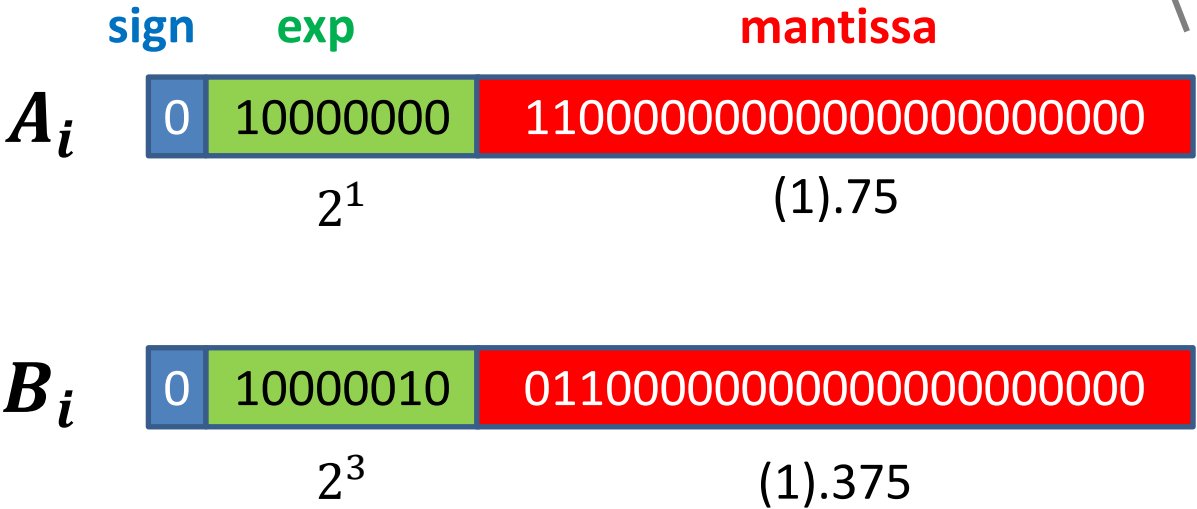
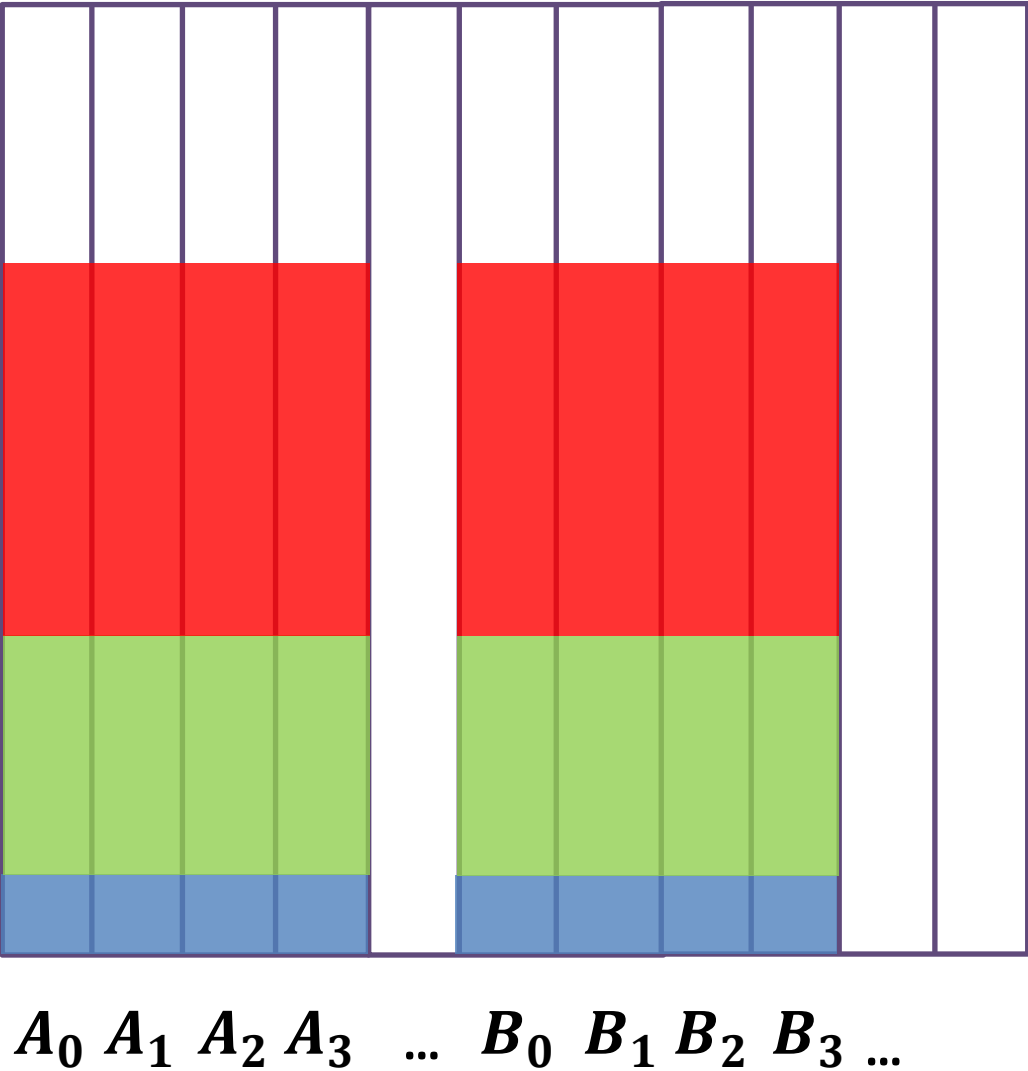
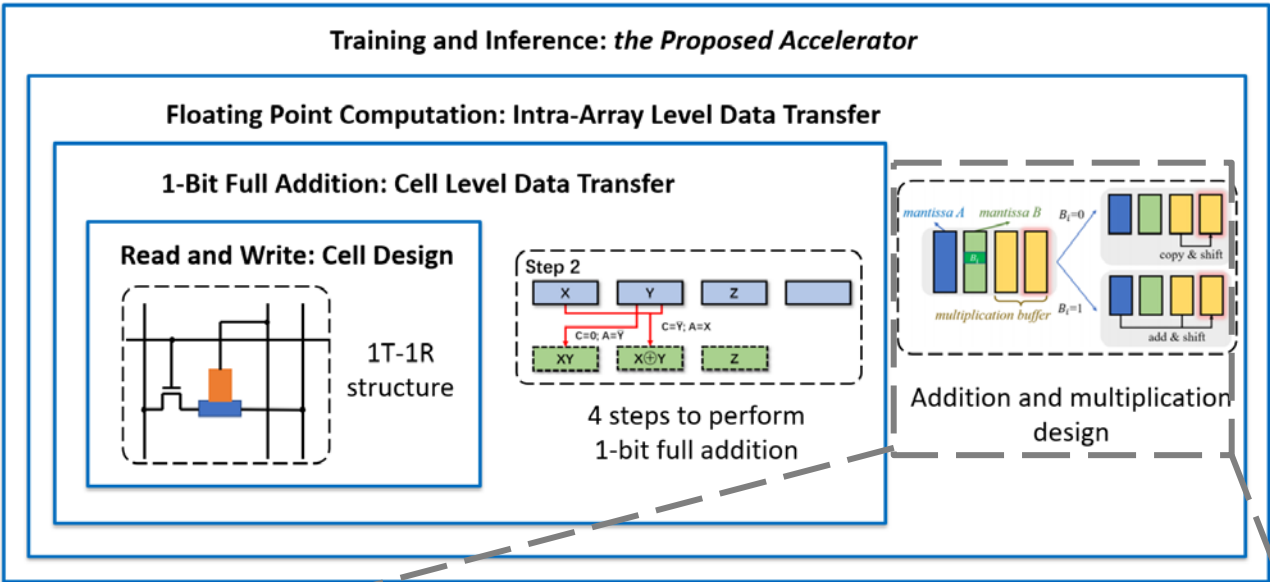
FloatPIM [Imani, ISCA'19]

$$S(R) = \overline{\bar{X} + \bar{Y} + \bar{Z}} + \overline{X + Y + Z} + \overline{\overline{X + Y} + X + Z} + \overline{\overline{Y + Z}}$$
$$Z' = \overline{\overline{X + Y} + X + Z} + \overline{\overline{Y + Z}}$$

The Proposed Method

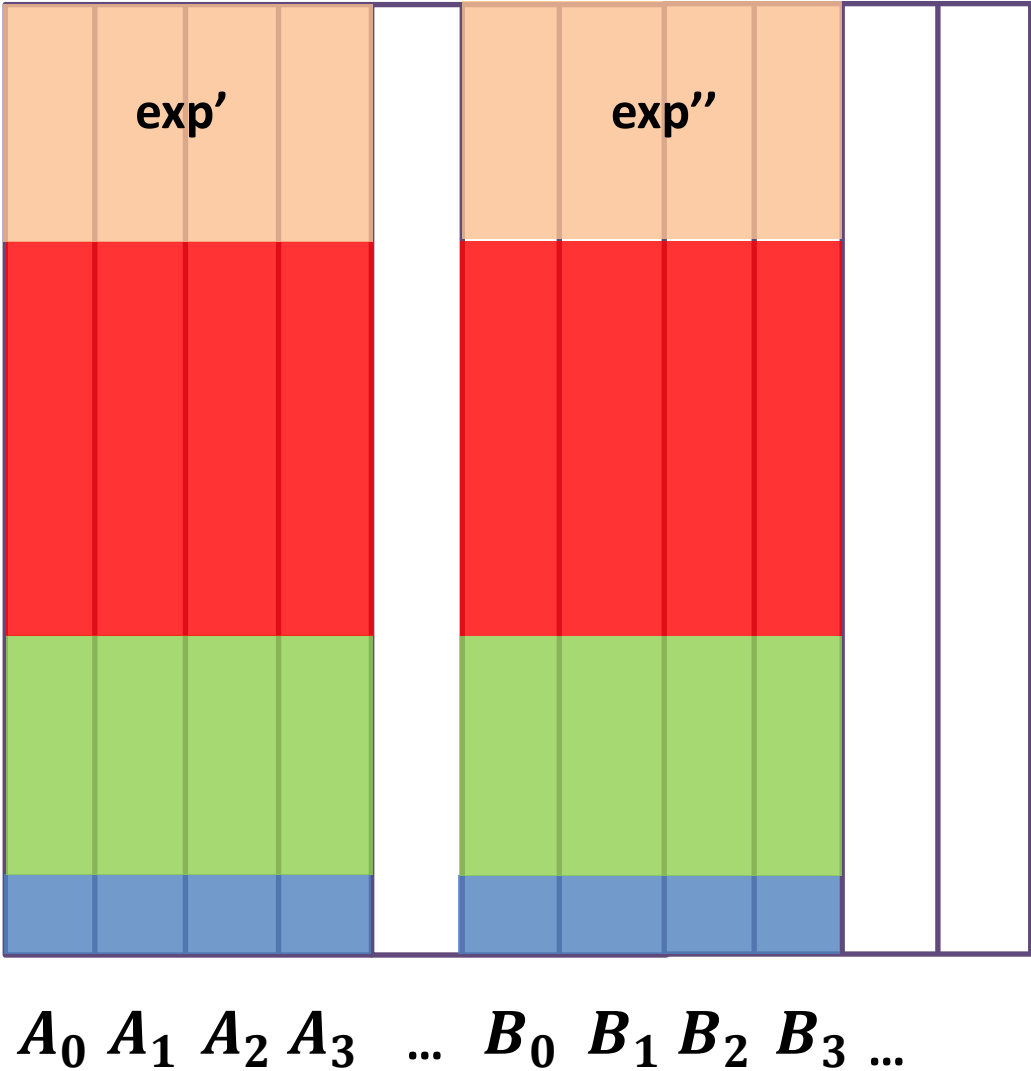
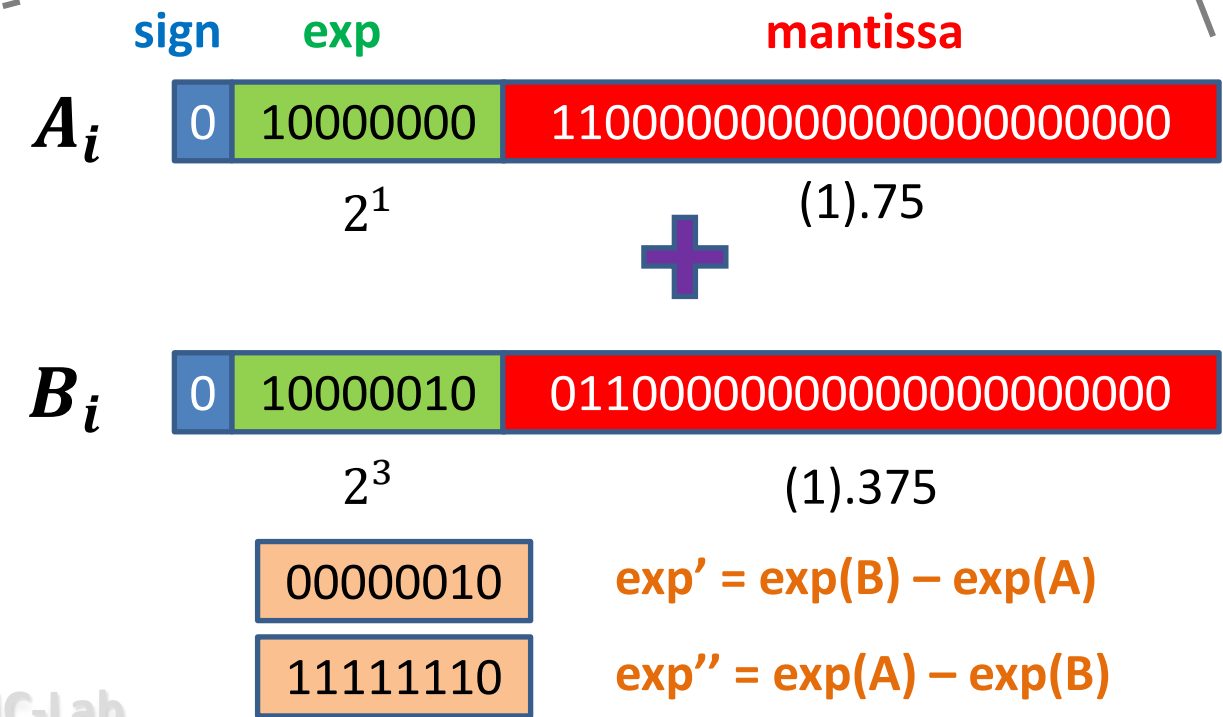
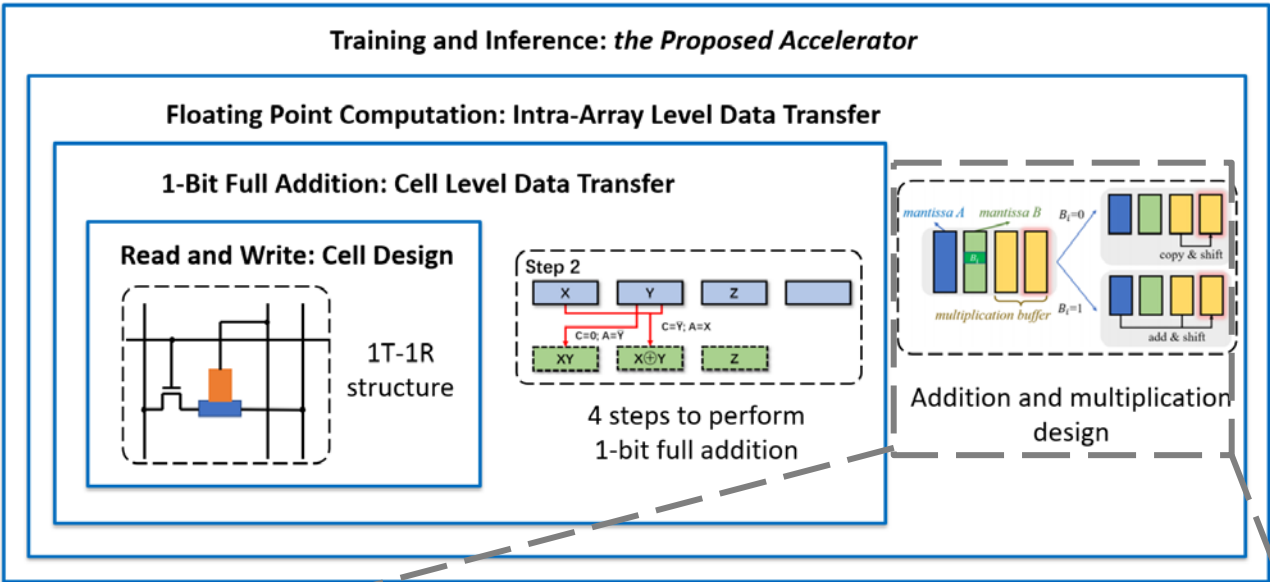
$$\begin{aligned} S(R) &= X \oplus Y \oplus Z \\ Z' &= XY + Z(X \oplus Y) \end{aligned}$$

Contribution 3: Improved Floating Point Computation



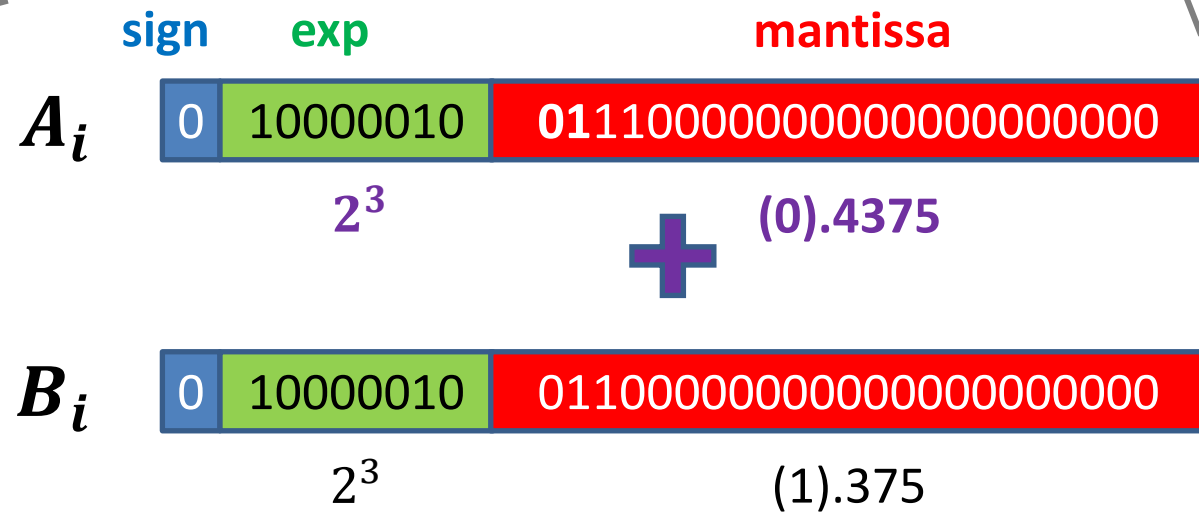
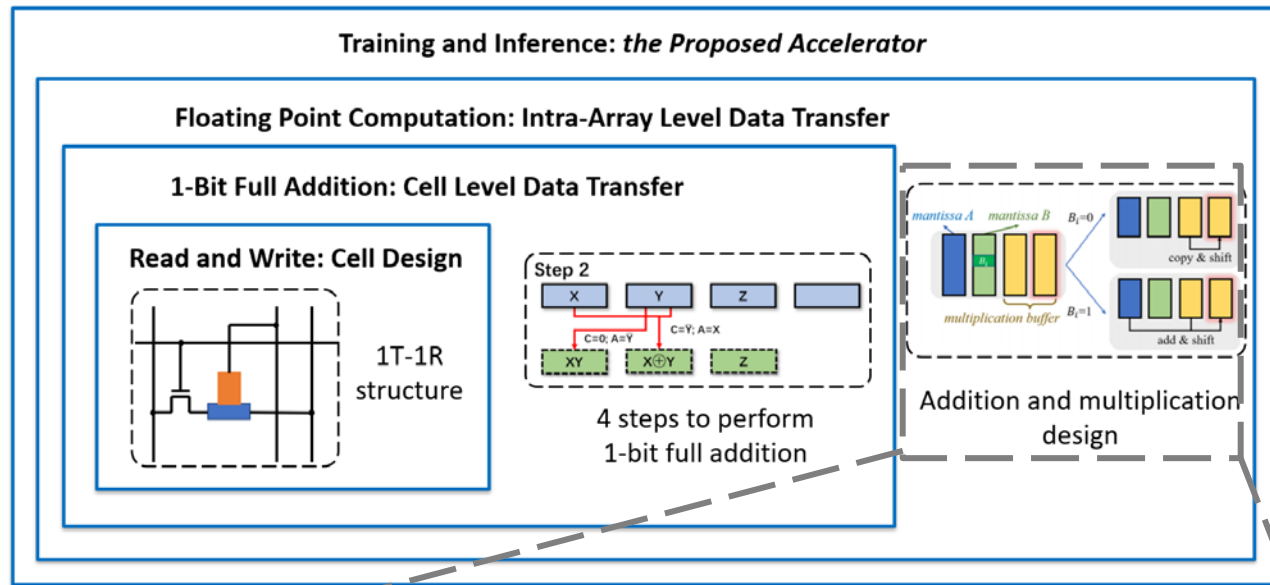
Floating Point Addition: $A_i + B_i$

Contribution 3: Improved Floating Point Computation

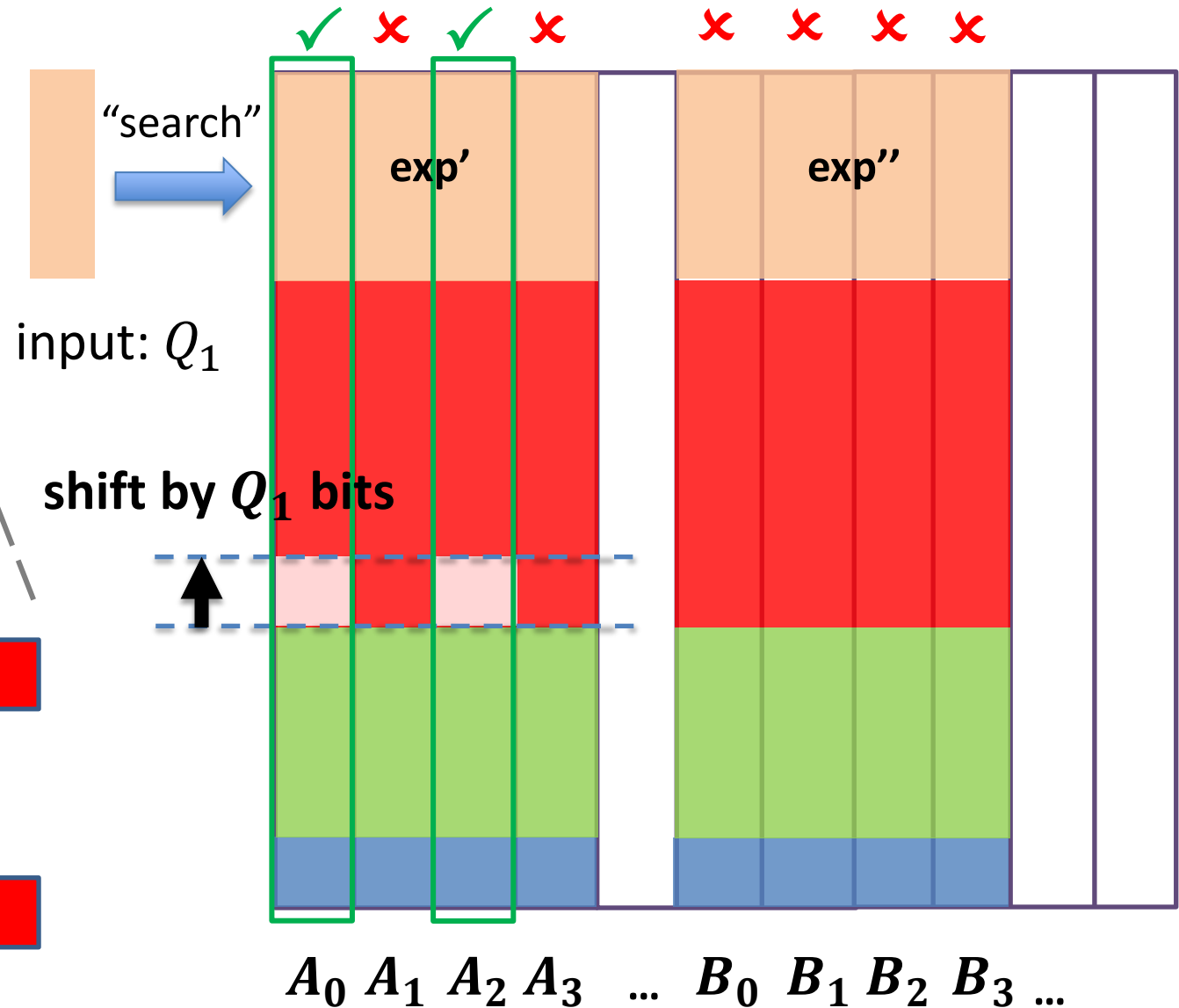


Column wise parallelism

Contribution 3: Improved Floating Point Computation

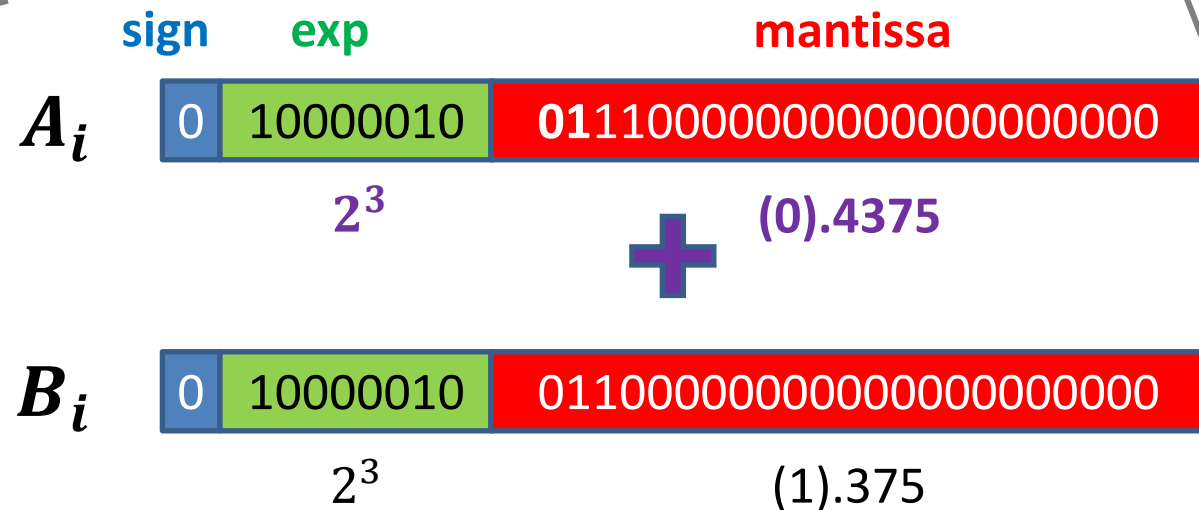
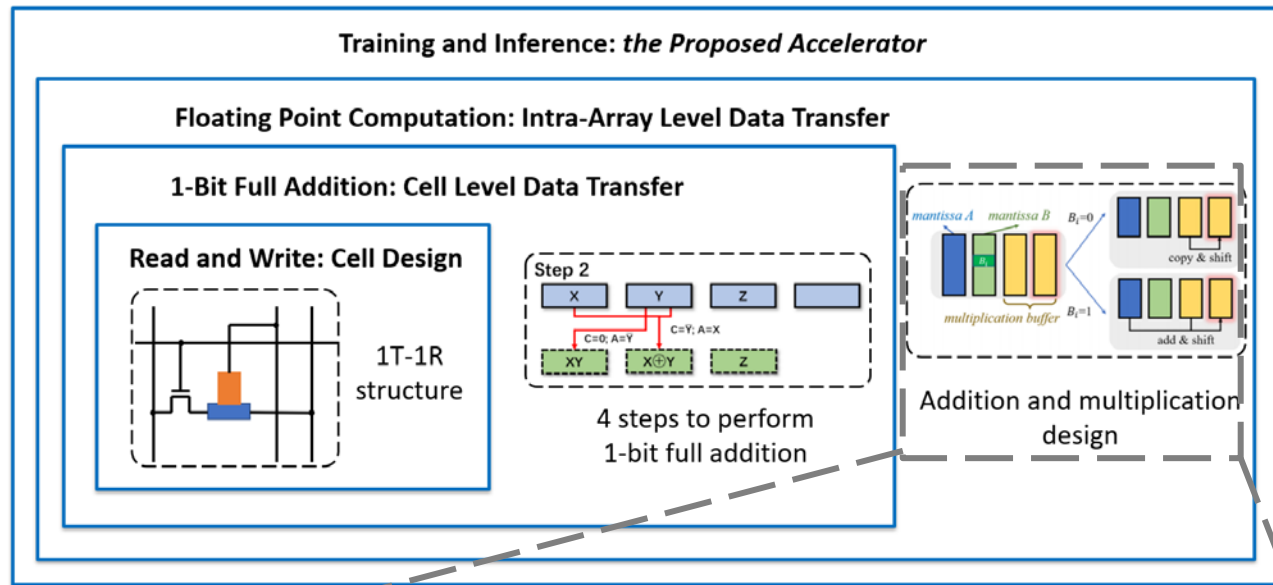


right shift mantissa to align exponents

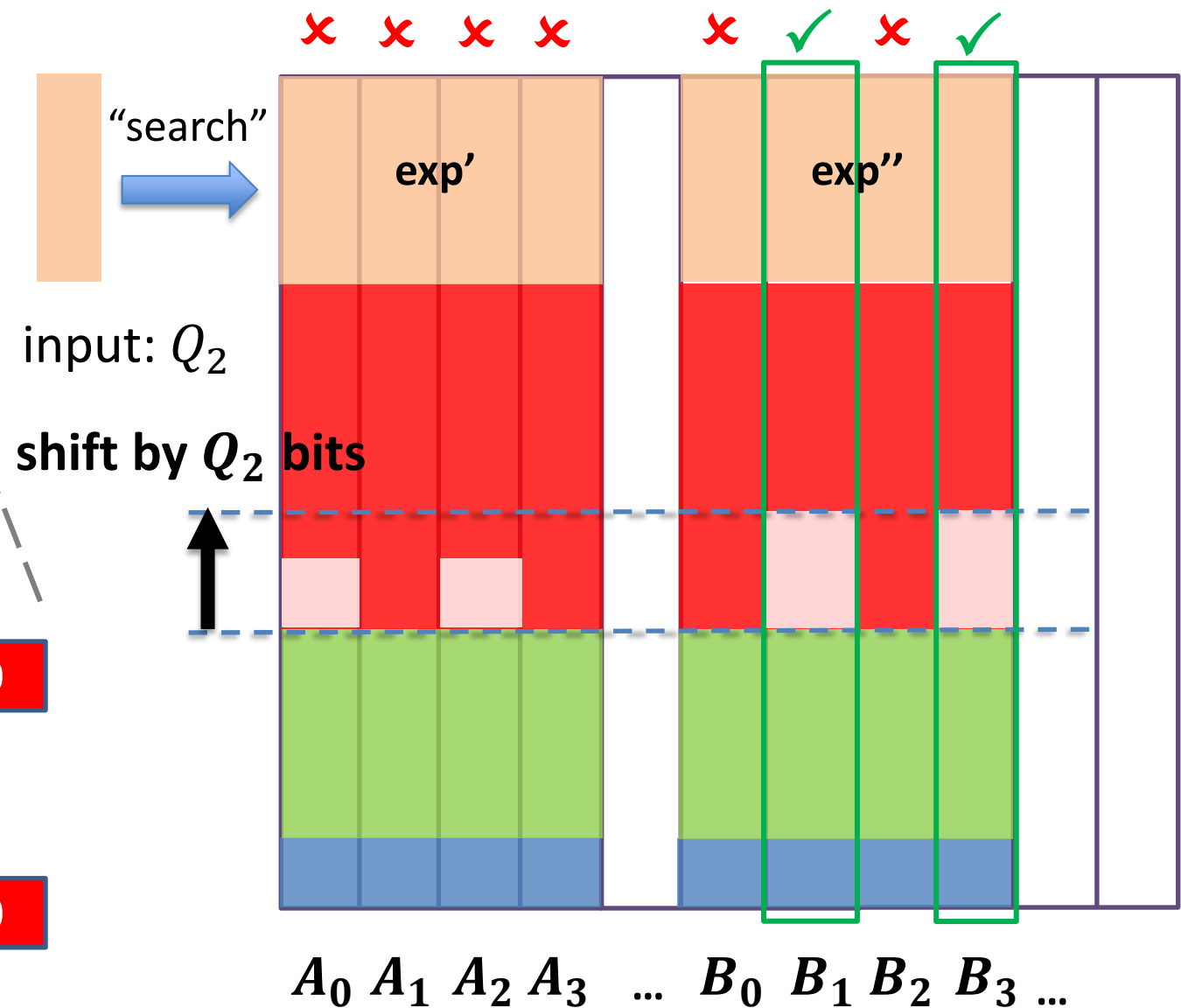


Column wise parallelism
(for same exp' or exp'')

Contribution 3: Improved Floating Point Computation

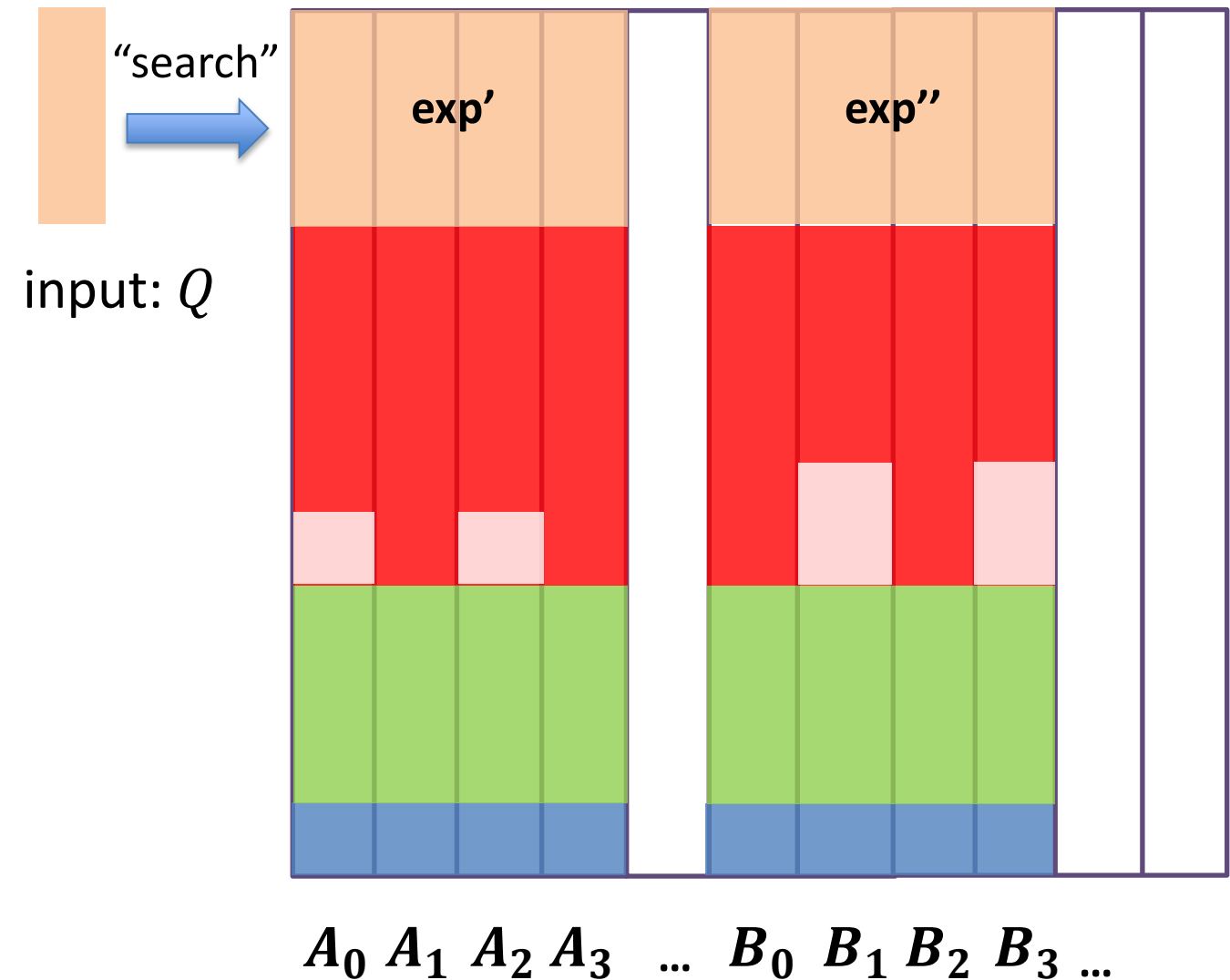
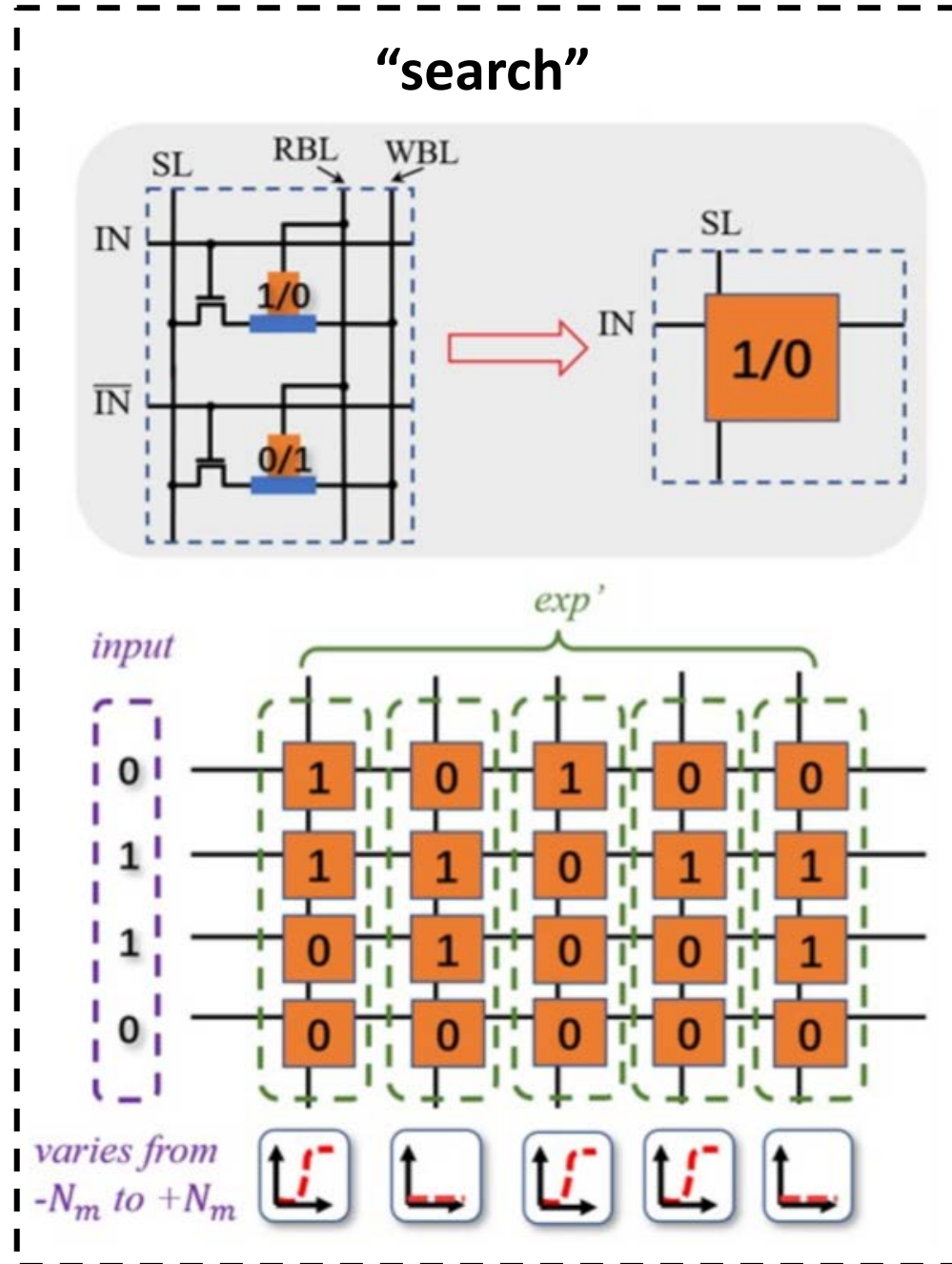


right shift mantissa to align exponents



Column wise parallelism
(for same exp' or exp'')

Contribution 3: Improved Floating Point Computation



Column wise parallelism
(for same exp' or exp'')

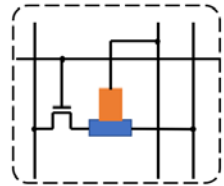
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

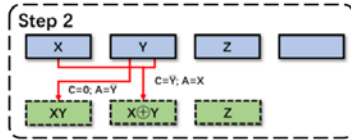
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

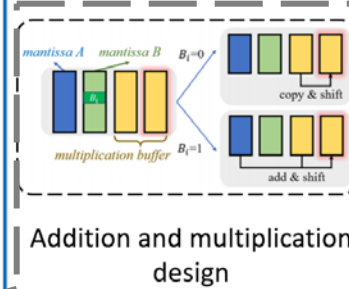
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Addition and multiplication
design

Floating Point Multiplication

XOR between sign bits

exponents addition

mantissas multiplication

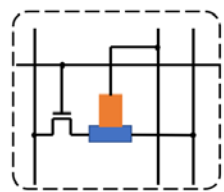
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

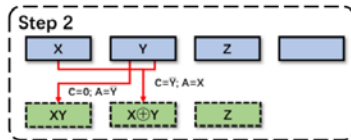
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

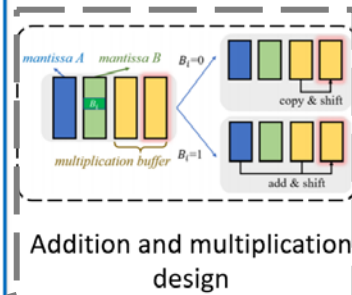
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Floating Point Multiplication

XOR between sign bits

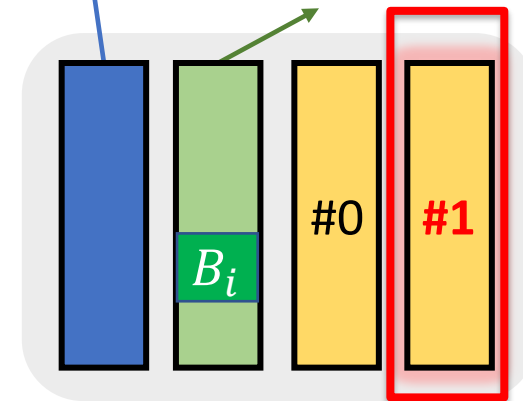
exponents addition

mantissas multiplication

Target Buffer in a Ping-Pong Manner

mantissa A

mantissa B



*buffer for
multiplication*

target buffer

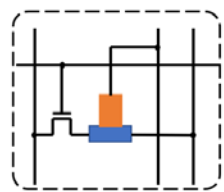
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

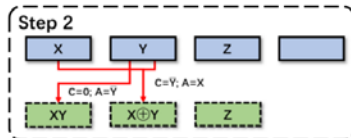
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

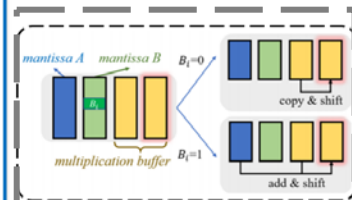
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Addition and multiplication
design

Floating Point Multiplication

XOR between sign bits

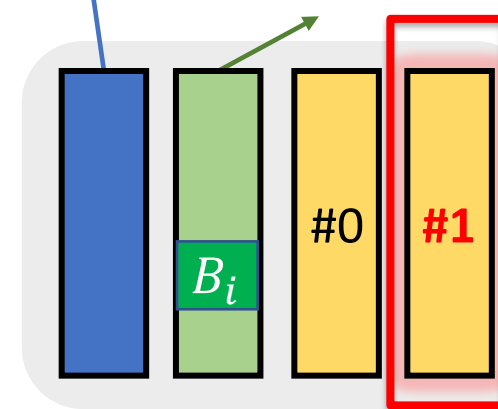
exponents addition

mantissas multiplication

Target Buffer in a Ping-Pong Manner

mantissa A

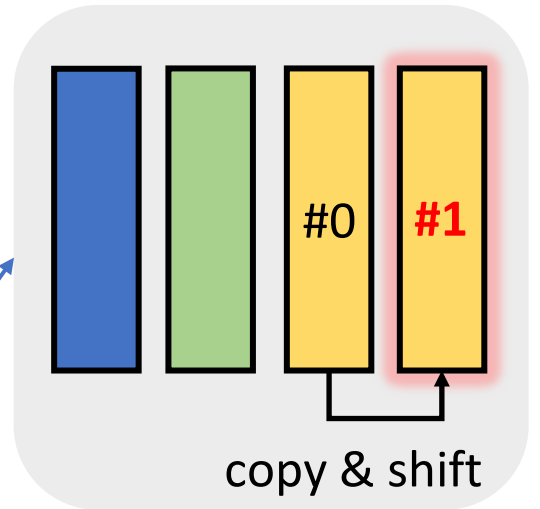
mantissa B



*buffer for
multiplication*

target buffer

$B_i=0$



copy & shift

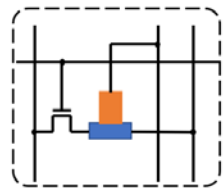
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

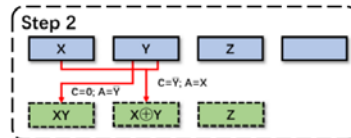
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

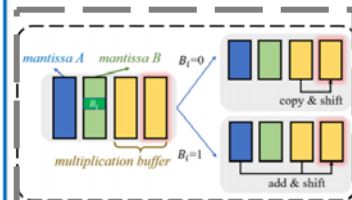
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Addition and multiplication
design

Floating Point Multiplication

XOR between sign bits

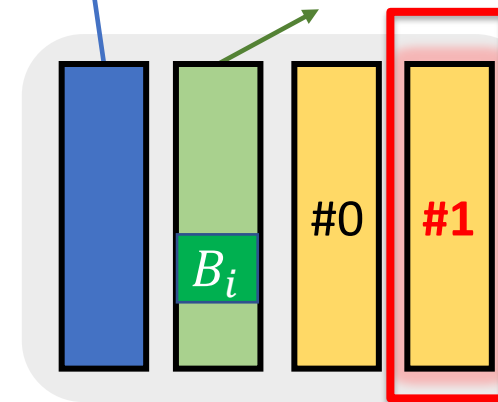
exponents addition

mantissas multiplication

Target Buffer in a Ping-Pong Manner

mantissa A

mantissa B

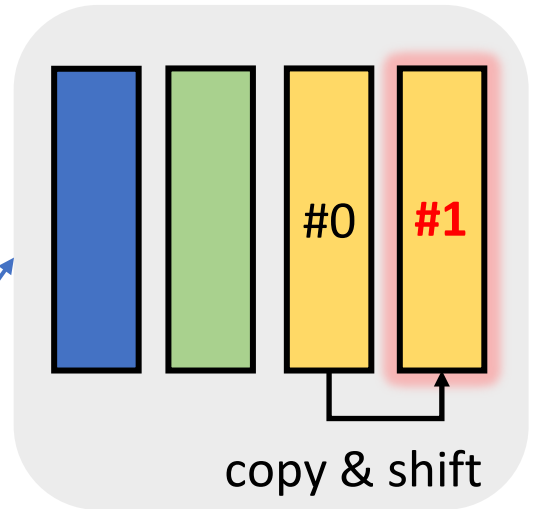


*buffer for
multiplication*

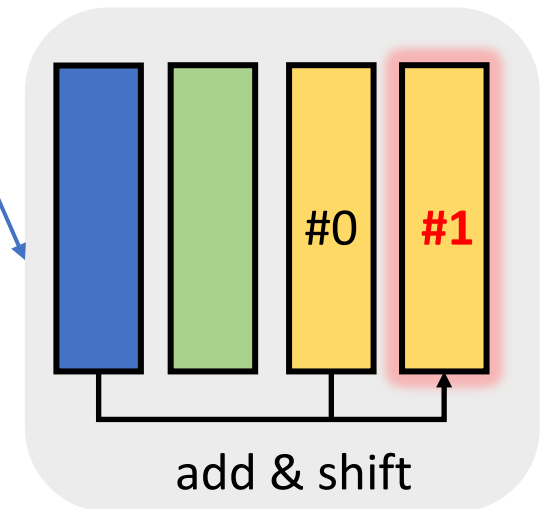
target buffer

$B_i=0$

$B_i=1$



copy & shift



add & shift

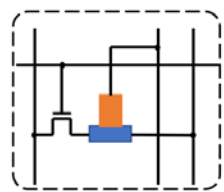
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

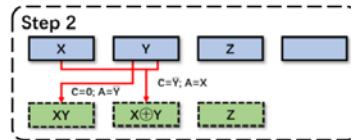
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

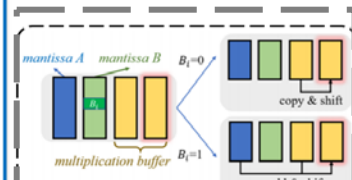
Read and Write: Cell Design



1T-1R structure



4 steps to perform
1-bit full addition



Addition and multiplication design

Floating Point Multiplication

XOR between sign bits

exponents addition

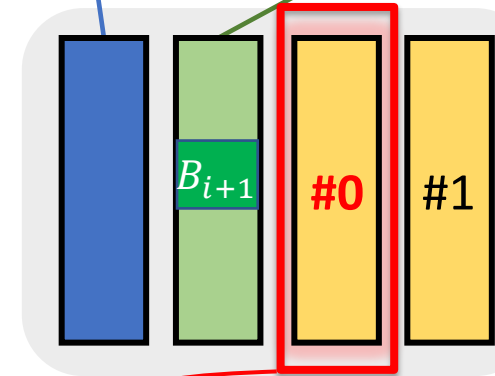
mantissas multiplication

Target Buffer in a Ping-Pong Manner

mantissa A

mantissa B

$B_{i+1}=0$



*buffer for
multiplication*

$B_{i+1}=1$

target buffer

copy & shift

add & shift

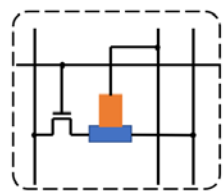
Contribution 3: Improved Floating Point Computation

Training and Inference: *the Proposed Accelerator*

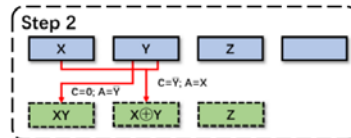
Floating Point Computation: Intra-Array Level Data Transfer

1-Bit Full Addition: Cell Level Data Transfer

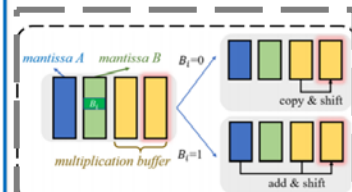
Read and Write: Cell Design



1T-1R
structure



4 steps to perform
1-bit full addition



Addition and multiplication
design

Floating Point Multiplication

XOR between sign bits

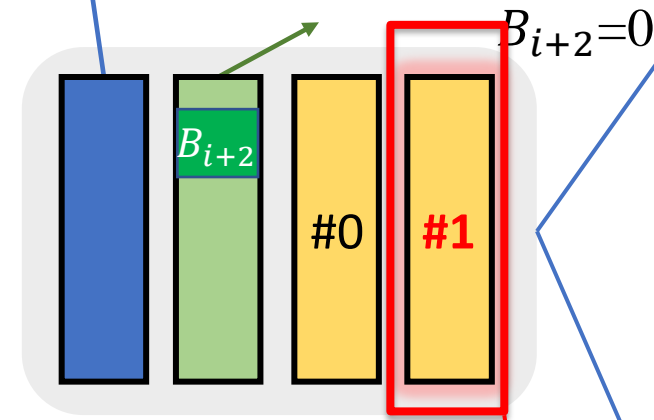
exponents addition

mantissas multiplication

Target Buffer in a Ping-Pong Manner

mantissa A

mantissa B



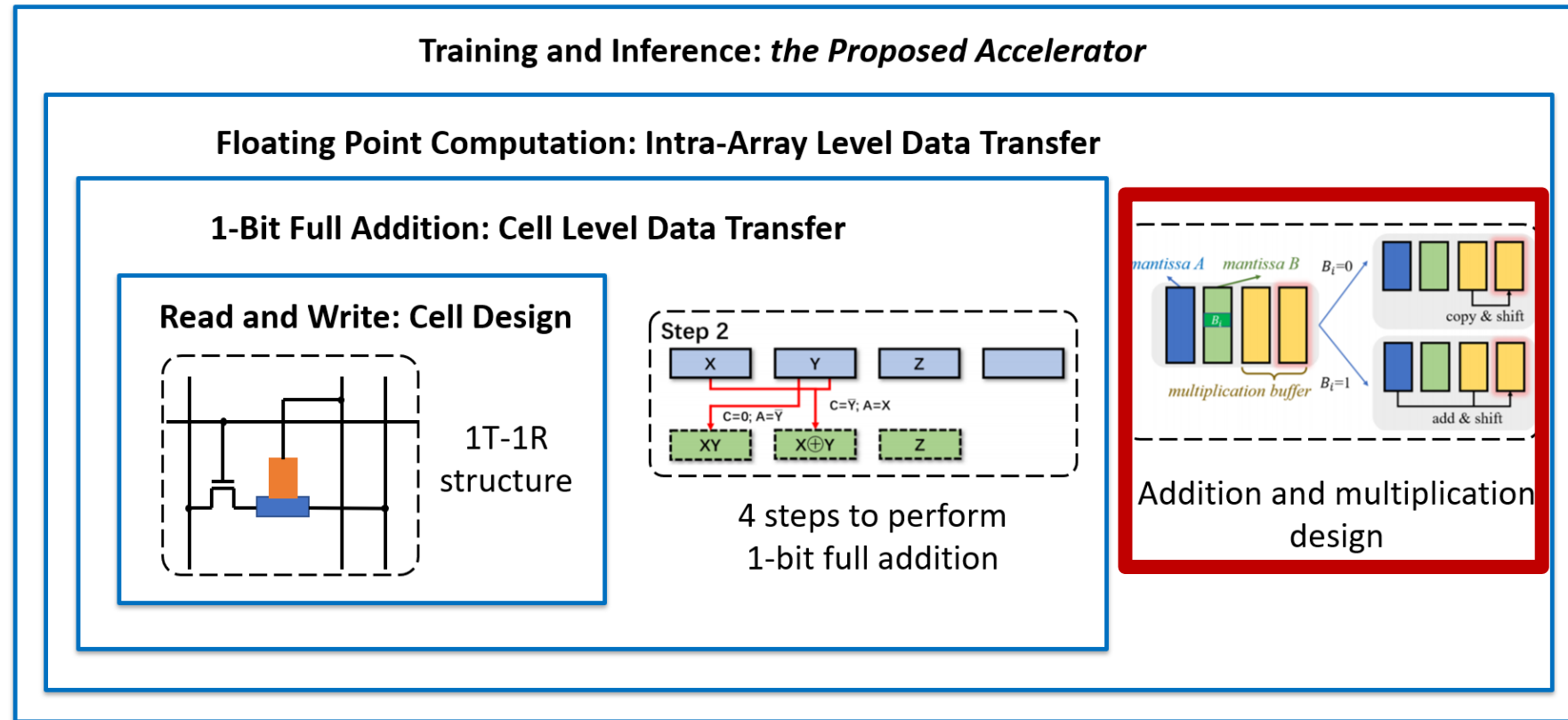
*buffer for
multiplication*

target buffer

copy & shift

add & shift

The Proposed Over SOTA Floating Point Computation



Compare with state-of-the-art work **FloatPIM** [Imani, ISCA'19]

- **Addition**: # of read and write clock cycles reduced by 80.5%↓, from $O(N^2)$ to $O(N)$
- **Multiplication**: # of read and write clock cycles reduced by 81.1%↓, # of memory cells requirement reduced by 77.4%↓

— for 32-bit floating point numbers

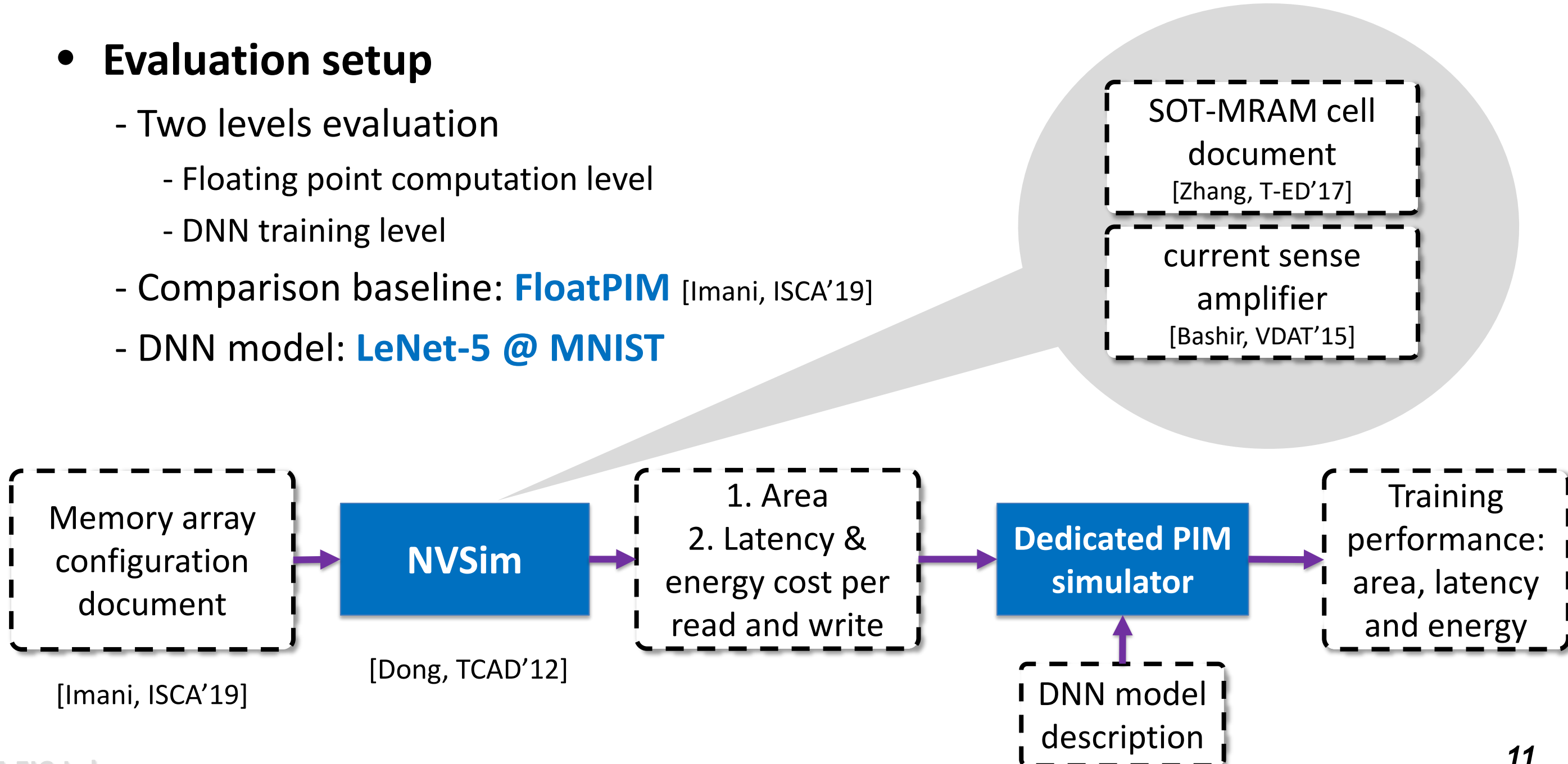
Outline

- Background and Motivation
- *The Proposed Accelerator: Overview*
- *The Proposed Accelerator: Computing Methodology*
- **Evaluation of *the Proposed Accelerator***
- Conclusion

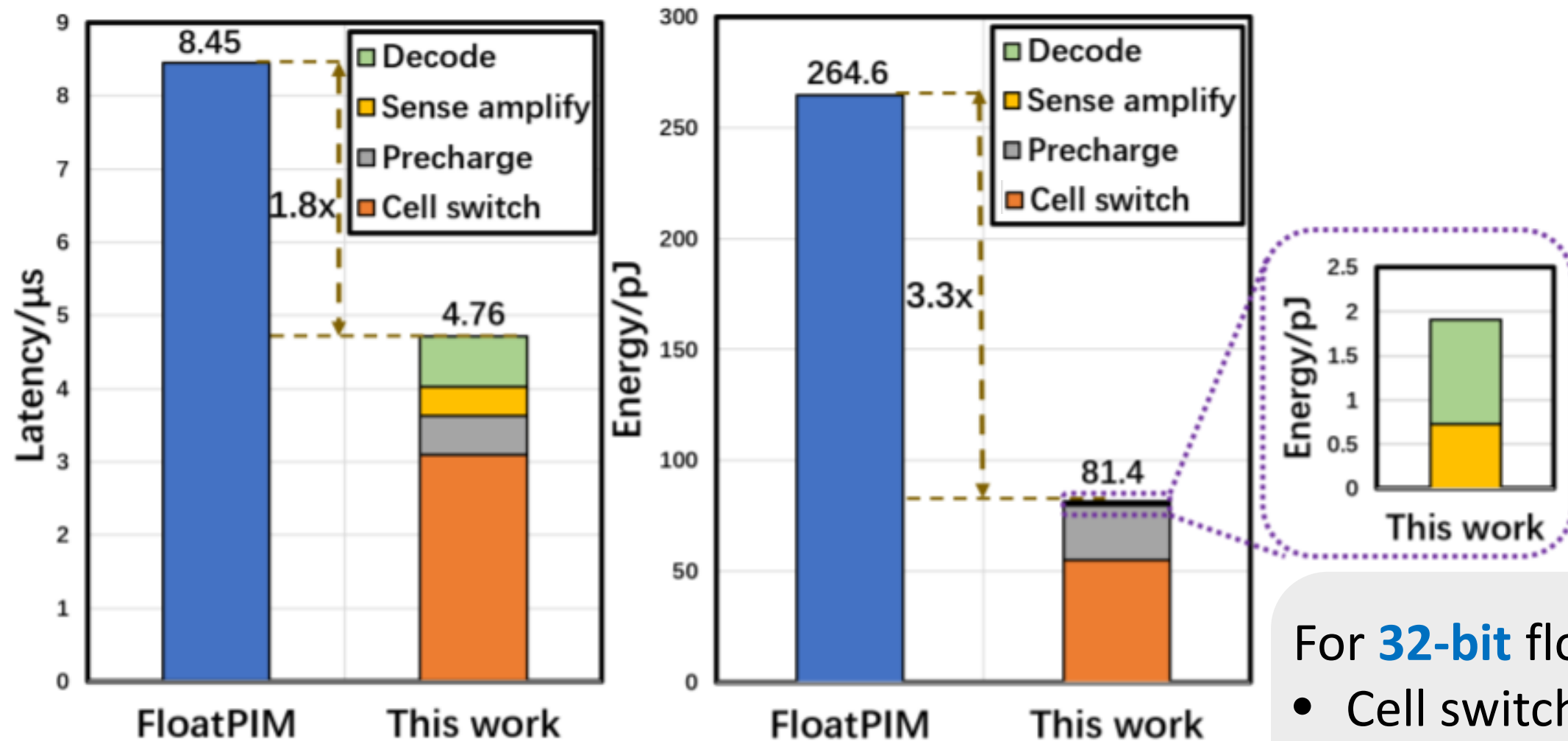
Evaluation Setting of *the Proposed Accelerator*

- **Evaluation setup**

- Two levels evaluation
 - Floating point computation level
 - DNN training level
- Comparison baseline: **FloatPIM** [Imani, ISCA'19]
- DNN model: **LeNet-5 @ MNIST**



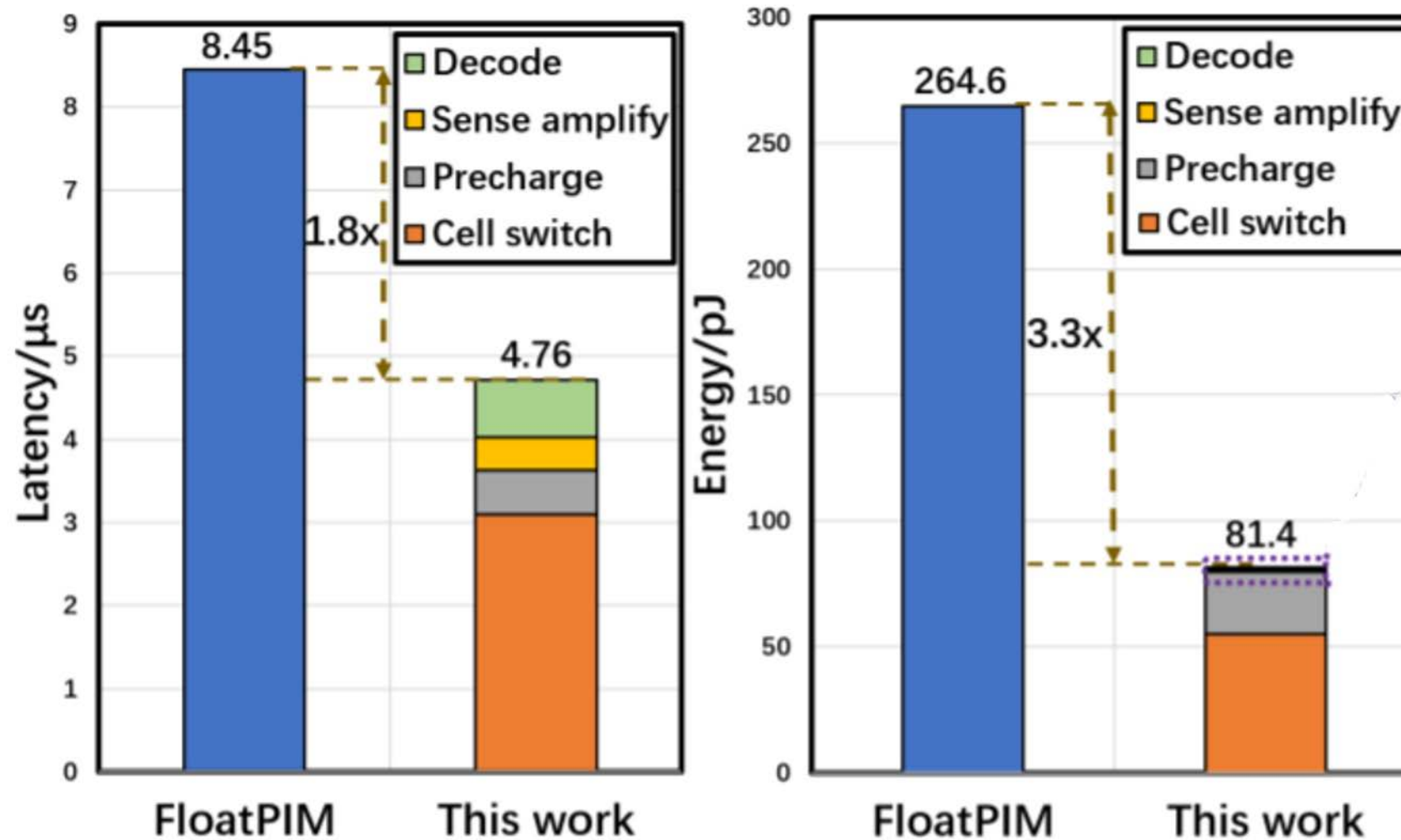
Evaluation Results on MAC



Multiplication and Accumulation Calculation (MAC)
performance comparison

- For **32-bit** floating point numbers:
- Cell switch dominates overhead
 - Latency efficiency: $\uparrow 1.8x$
 - Energy efficiency: $\uparrow 3.3x$

Evaluation Results on MAC



Benefits from:

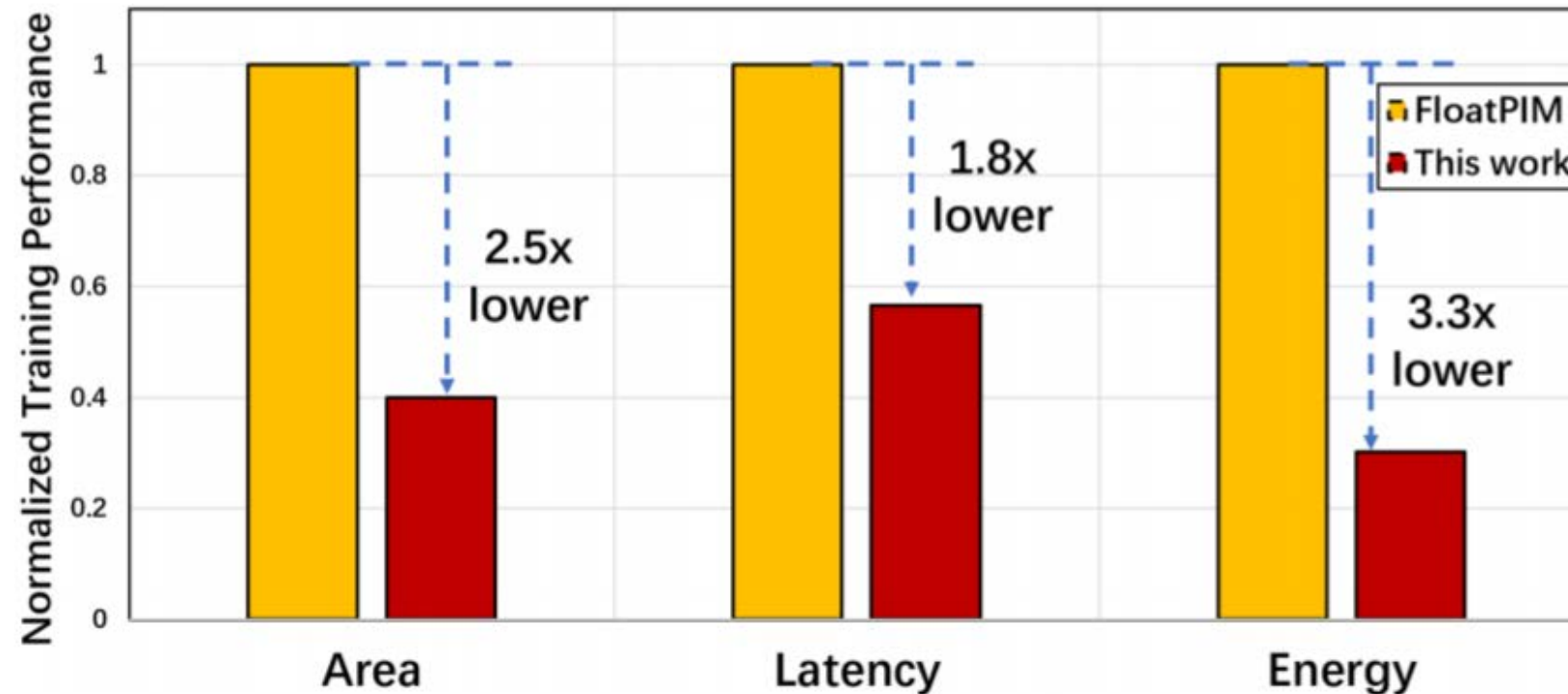
- **Fewer clock cycles** required due to improved full addition and floating point computation design
- **Low write current** due to SOT-MRAM based memory cell design

For **32-bit** floating point numbers:

- Cell switch dominates overhead
- Latency efficiency: **↑1.8x**
- Energy efficiency: **↑3.3x**

Multiplication and Accumulation Calculation (MAC)
performance comparison

Evaluation Results on Neural Network Training



Benefits from:

- **Fewer extra cells required** to store intermediate results due to improved full addition and floating point computation design
- **High flexibility** memory assignment to **maximize cells reuse** due to SOT-MRAM based 1T-1R cell design
- The improvement of MAC

LeNet-5 @ MNIST

- Area efficiency: **↑2.5x**
- Latency efficiency: **↑1.8x**
- Energy efficiency: **↑3.3x**

Outline

- Background and Motivation
- *The Proposed Accelerator: Overview*
- *The Proposed Accelerator: Computing Methodology*
- Evaluation of *the Proposed Accelerator*
- Conclusion

Conclusion

- ***The Proposed Accelerator***: A new SOT-MRAM based Process In-Memory accelerator for efficient training with floating point precision
 - New memory cell design which features **an improved balance** between computation flexibility and memory density
 - New full addition design that **improves computation efficiency**
 - New floating point computation design for **efficient training**
- Achieve **3.3×, 1.8×, and 2.5×** improvement in terms of energy, latency, and area, respectively, compared with a state-of-the-art PIM based DNN training accelerator

MRAM based PIM Accelerator for NN Training

*Hongjie Wang**, *Yang Zhao**, Chaojian Li, Yue Wang, and **Yingyan Lin**



Questions ?



*the first two authors contribute equally