



**JULY 9-13, 2023**

**MOSCONE WEST CENTER  
SAN FRANCISCO, CA, USA**







# *BP-NTT: Fast and Compact in-SRAM Number Theoretic Transform with Bit-Parallel Modular Multiplication*

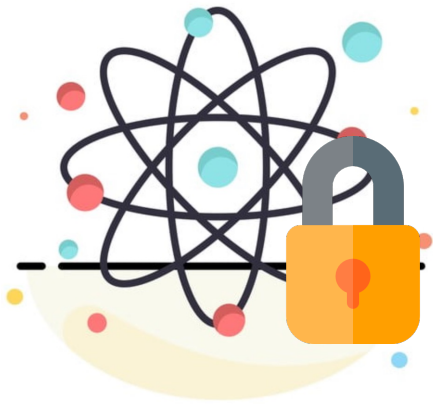
Jingyao Zhang<sup>\*</sup>, Mohsen Imani<sup>†</sup>, Elaheh Sadredini<sup>\*</sup>



# Lattice-based Cryptography is the Future

---

**NIST**



Post-quantum Cryptography



Homomorphic Encryption

# Lattice-based Cryptography is the Future

---

- **Polynomial Multiplication** is the bedrock

**NIST**



Post-quantum Cryptography



Homomorphic Encryption

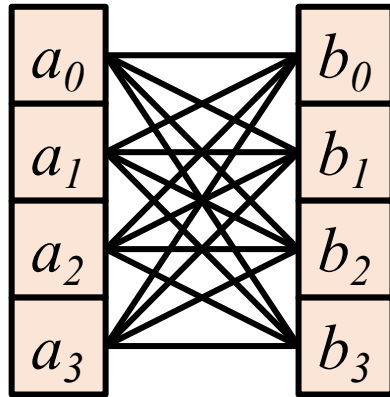
**Polynomial Multiplication**

# Polynomial Multiplication is the Bedrock

---

# Polynomial Multiplication is the Bedrock

---

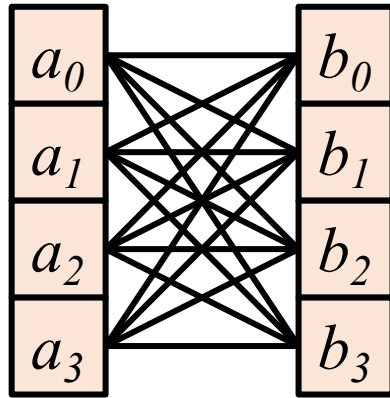


Polynomial Multiplication

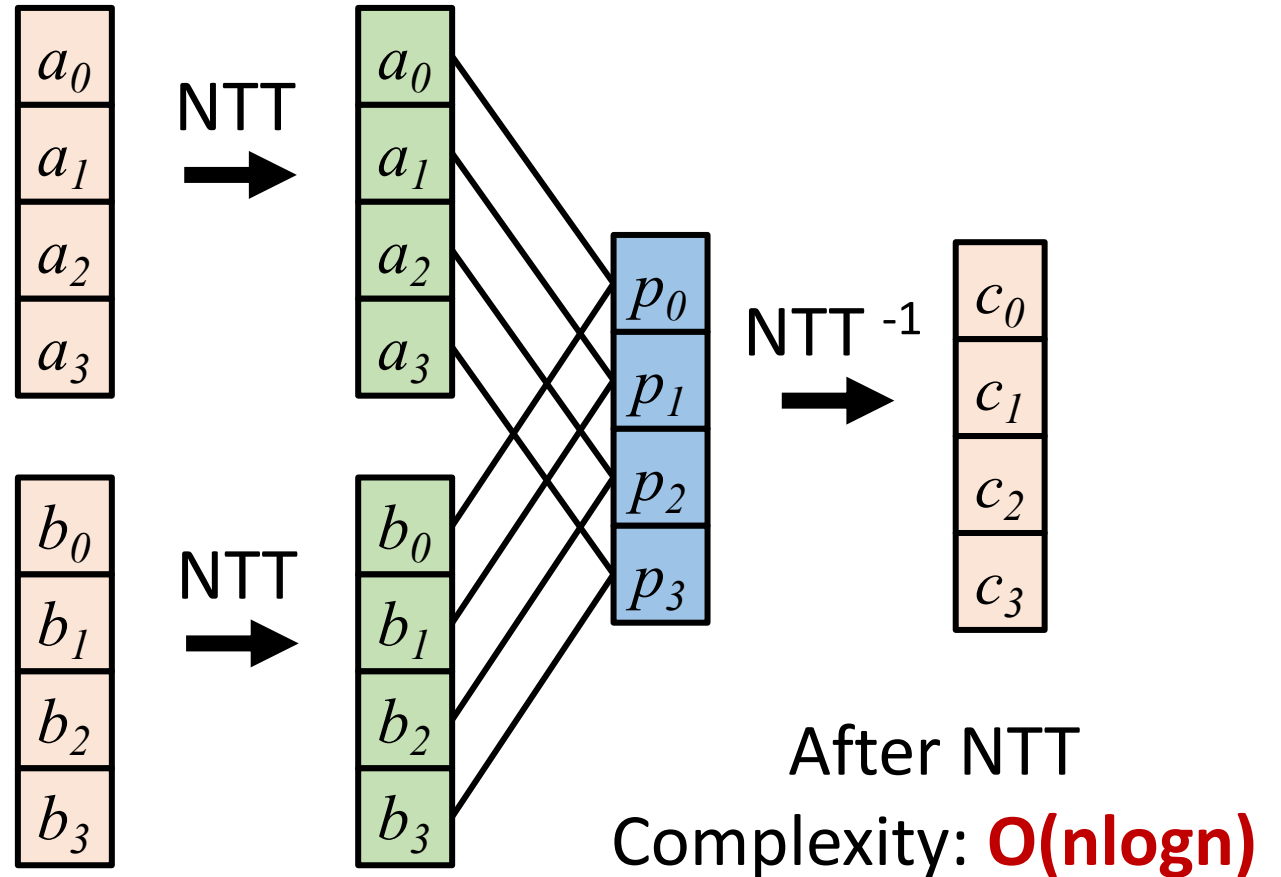
Complexity:  $O(n^2)$

# Polynomial Multiplication is the Bedrock

- Number Theoretic Transform (NTT) is necessary



Polynomial Multiplication  
Complexity:  $O(n^2)$



# NTT Acceleration is Essential

---



# NTT Acceleration is Essential

---

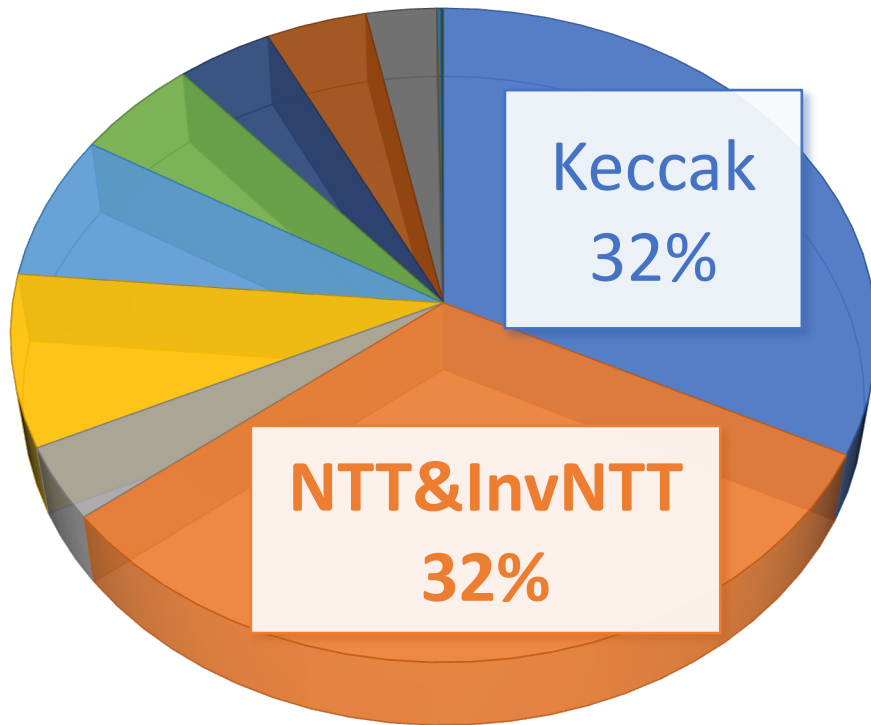
- From profiling, 32~50% of execution time is spent on NTT/InvNTT

# NTT Acceleration is Essential

---

- From profiling, 32~50% of execution time is spent on NTT/InvNTT

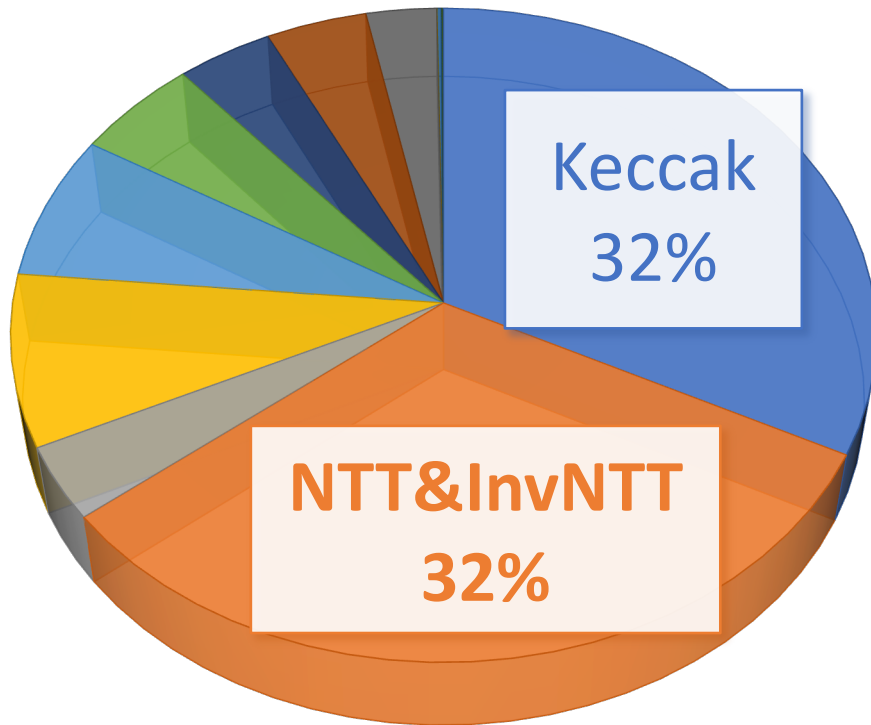
CRYSTAL-KYBER



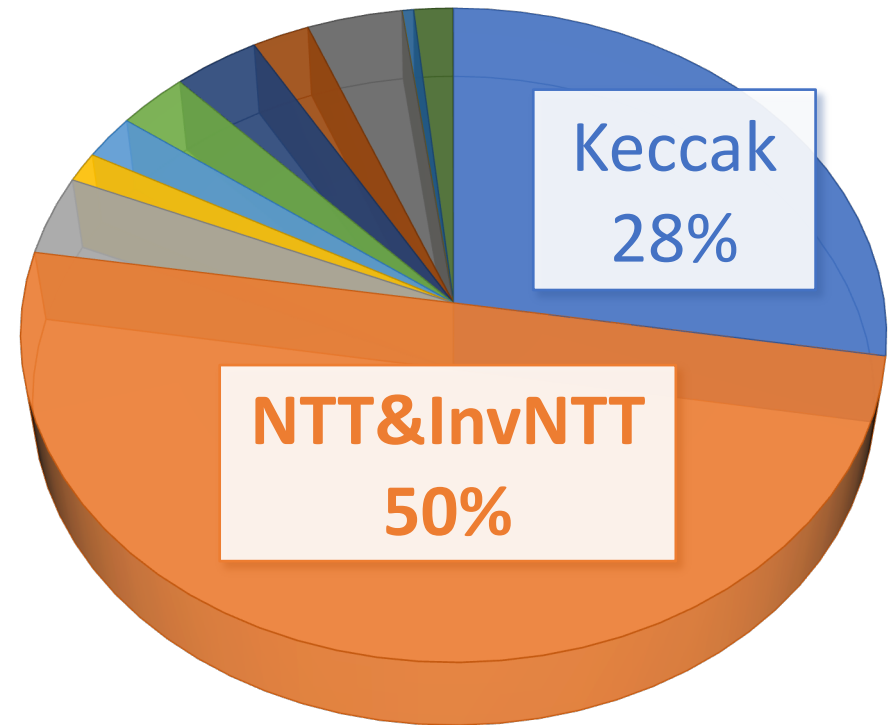
# NTT Acceleration is Essential

- From profiling, 32~50% of execution time is spent on NTT/InvNTT

CRYSTAL-KYBER

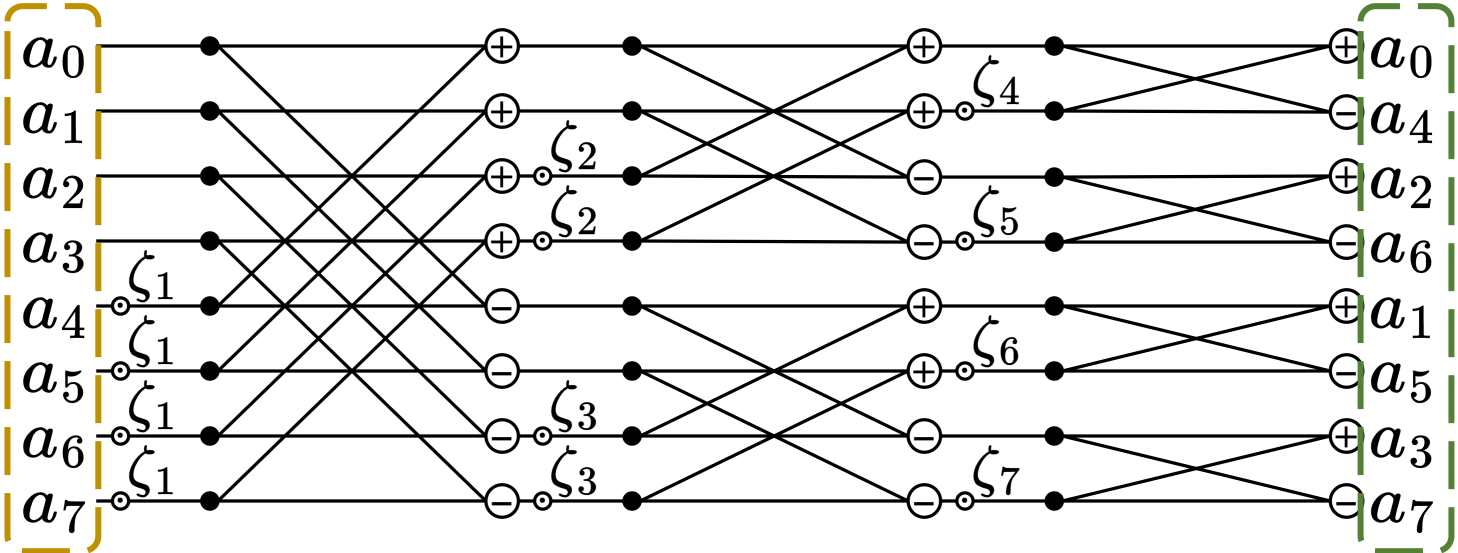


CRYSTAL-DILITHIUM



# NTT Acceleration is **Hard!**

Polynomial coefficients before NTT



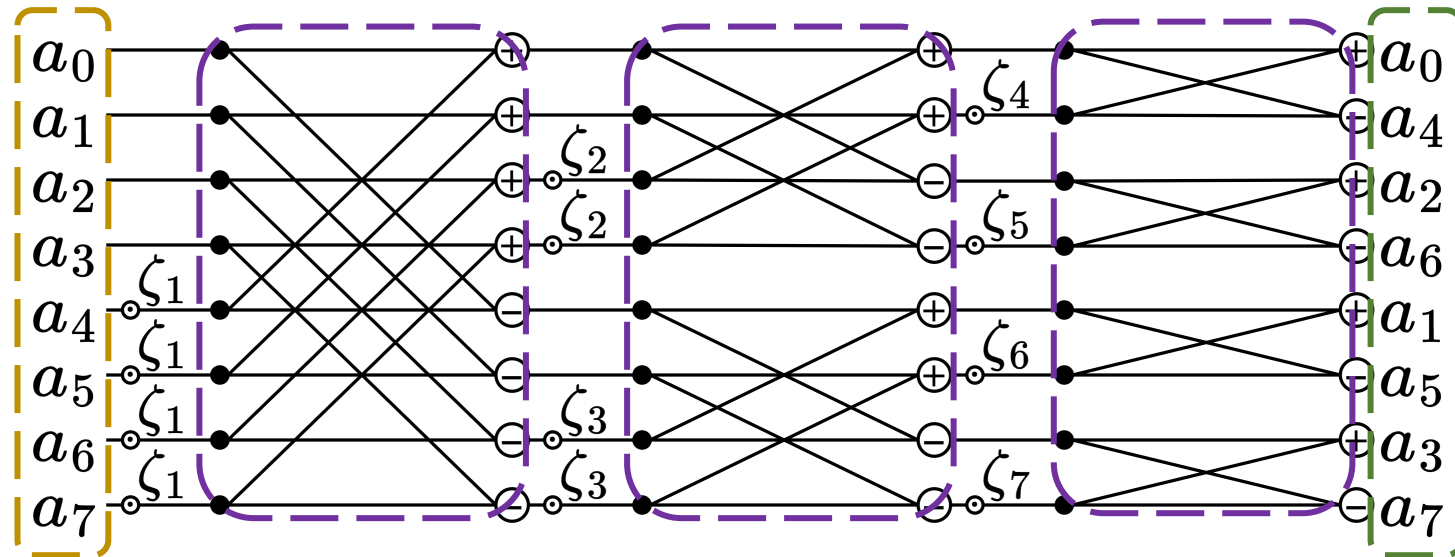
Polynomial coefficients after NTT

*3-stage Cooley-Tukey butterfly*

# NTT Acceleration is **Hard!**

- Complicated data dependencies

Polynomial  
coefficients  
before NTT



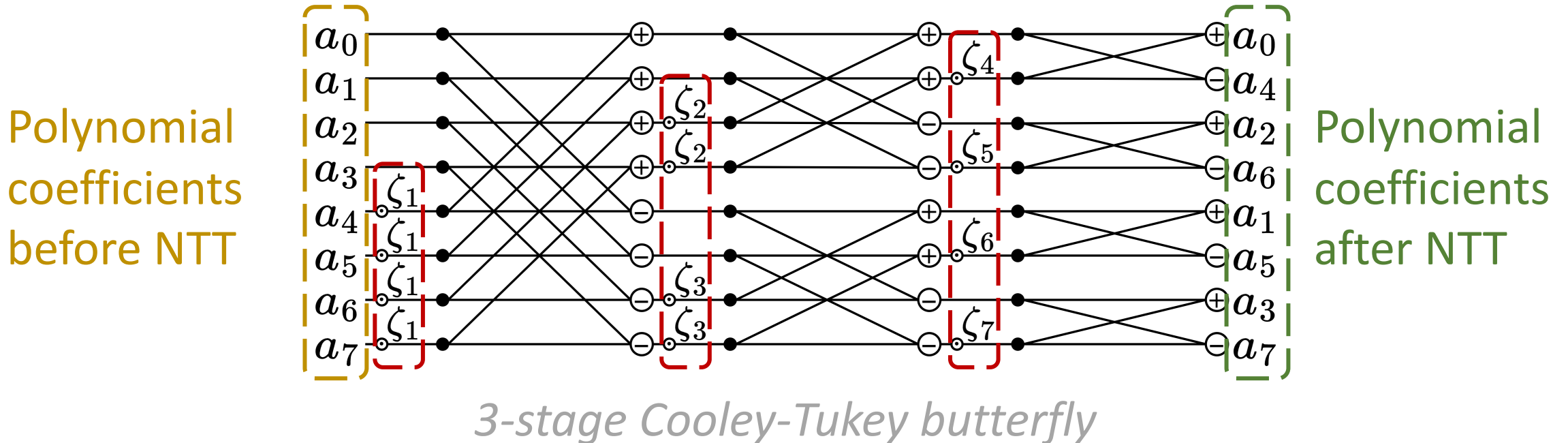
Polynomial  
coefficients  
after NTT

*3-stage Cooley-Tukey butterfly*



# NTT Acceleration is **Hard!**

- Complicated data dependencies
- Heavy multiplication with division-based modulo operation



# Existing Solutions

---

# Existing Solutions

---



## *ASIC*

LEIA [CICC '18], Sapphire [ISSCC '19], ...

# Existing Solutions

---



## *ASIC*

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement

# Existing Solutions

---



## ASIC

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement
- ❑ Limited flexibility
  - Dedicated modular multiplier



# Existing Solutions

---



## ASIC

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement
- ❑ Limited flexibility
  - Dedicated modular multiplier



How to achieve **efficient & flexible** NTT acceleration?

# Existing Solutions

---



## ASIC

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement
- ❑ Limited flexibility
  - Dedicated modular multiplier



## Compute-in-Memory

Recryptor [JSSC '18], Duality Cache [ISCA '19], ...



How to achieve **efficient & flexible** NTT acceleration?

# Existing Solutions

---



## ASIC

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement
- ❑ Limited flexibility
  - Dedicated modular multiplier



## Compute-in-Memory

Recryptor [JSSC '18], Duality Cache [ISCA '19], ...

- ❑ Potential high efficiency
  - Reduced data movement



How to achieve **efficient & flexible** NTT acceleration?

# Existing Solutions



## ASIC

LEIA [CICC '18], Sapphire [ISSCC '19], ...

- ❑ Low efficiency
  - Frequent data movement
- ❑ Limited flexibility
  - Dedicated modular multiplier



## Compute-in-Memory

Recryptor [JSSC '18], Duality Cache [ISCA '19], ...

- ❑ Potential high efficiency
  - Reduced data movement
- ❑ High flexibility
  - General vector processing units



How to achieve **efficient & flexible** NTT acceleration?

# Challenges for Memory-Centric NTT

---



# Challenges for Memory-Centric NTT

---

## Observations:

- ❑ Inefficient data layout incurs *redundant shifts*

# Challenges for Memory-Centric NTT

---

## Observations:

- ❑ Inefficient data layout incurs *redundant shifts*

## Contribution 1:

**Shift-optimized data layout to avoid *redundant shifts***

# Challenges for Memory-Centric NTT

---

## Observations:

- ❑ Inefficient data layout incurs *redundant shifts*

## Contribution 1:

**Shift-optimized data layout to avoid *redundant shifts***

- ❑ Bit-parallel modular multiplication has *unnecessary carry propagations*

# Challenges for Memory-Centric NTT

---

## Observations:

- Inefficient data layout incurs *redundant shifts*

## Contribution 1:

**Shift-optimized data layout** to avoid *redundant shifts*

- Bit-parallel modular multiplication has *unnecessary carry propagations*

## Contribution 2:

**Carry-save modular multiplication** to avoid *carry propagations*

# Overview of Our Solution: *BP-NTT*

---

# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

# Overview of Our Solution: *BP-NTT*

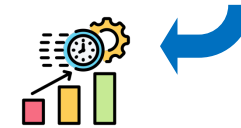
---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security



High Throughput



# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security  

High Throughput  

*BP-NTT* uses **shift-optimized data alignment**



# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security  

High Throughput  

*BP-NTT* uses **shift-optimized data alignment**

Low Latency  

  Low Energy

# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security  

High Throughput  

*BP-NTT* uses **shift-optimized data alignment**

Low Latency  

  Low Energy

*BP-NTT* employs **bit-parallel modular multiplication**

# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security  

High Throughput  

*BP-NTT* uses **shift-optimized data alignment**

Low Latency  

  Low Energy

*BP-NTT* employs **bit-parallel modular multiplication**

Low Latency  

  Small Area

# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes **LLC** to perform **bitline computing**

High security  

High Throughput  

*BP-NTT* uses **shift-optimized data alignment**

Low Latency  

 Low Energy 

*BP-NTT* employs **bit-parallel modular multiplication**

Low Latency  

 Small Area 

# *BP-NTT*: Repurposed LLC

---

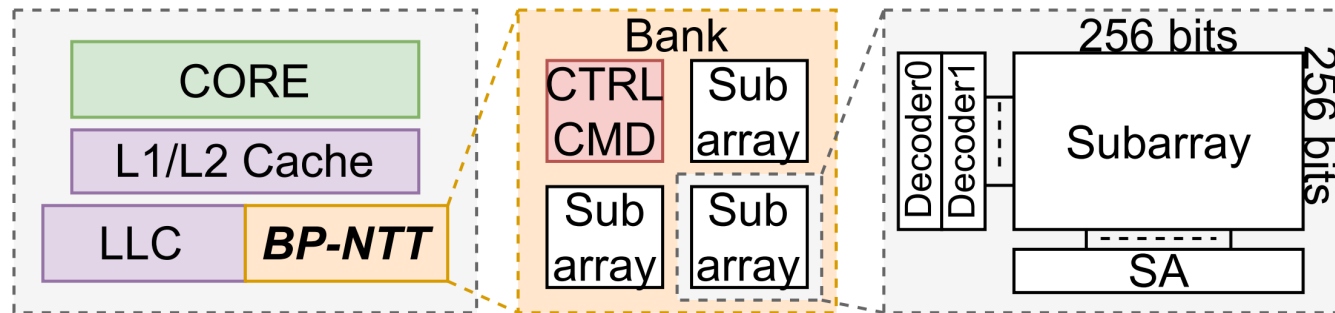
# *BP-NTT*: Repurposed LLC

---

- *BP-NTT* repurposes LLC to perform bitline computing [1]

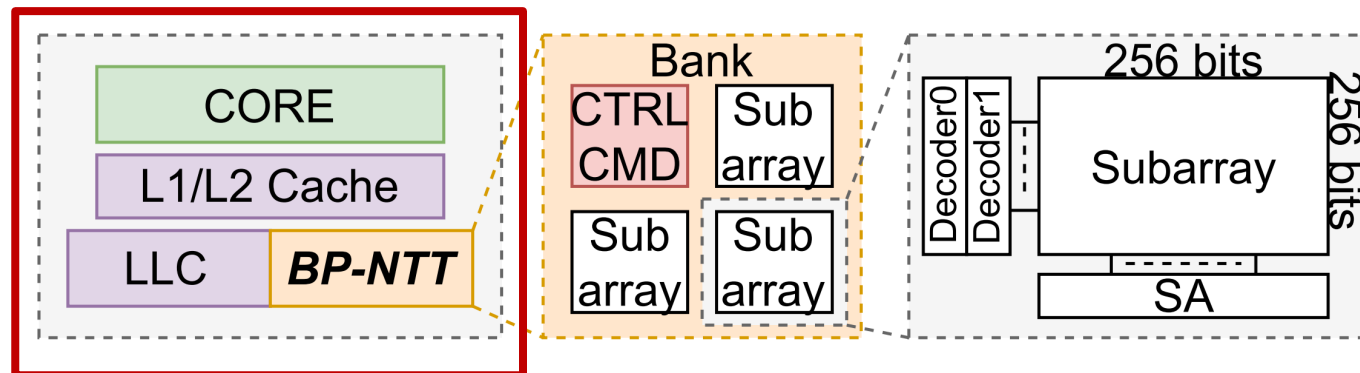
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]



# BP-NTT: Repurposed LLC

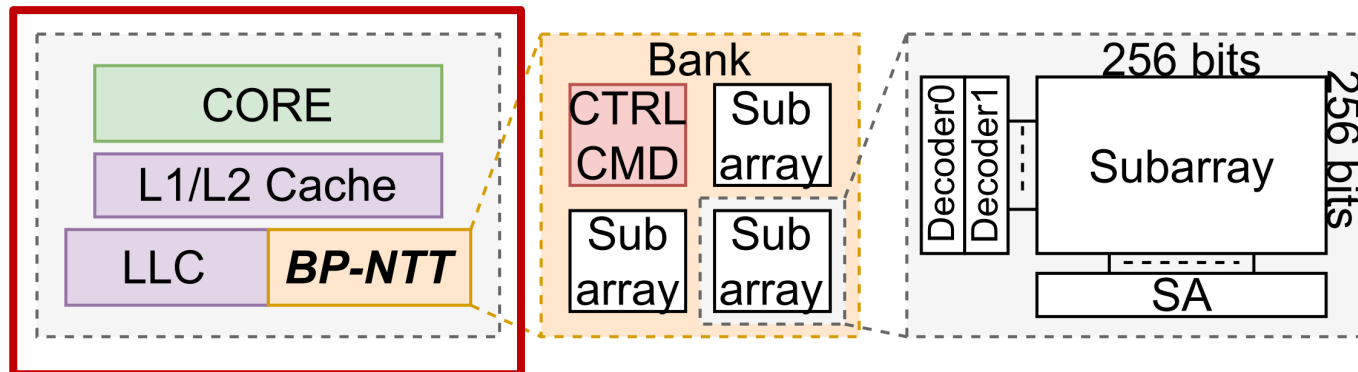
- *BP-NTT* repurposes LLC to perform bitline computing [1]





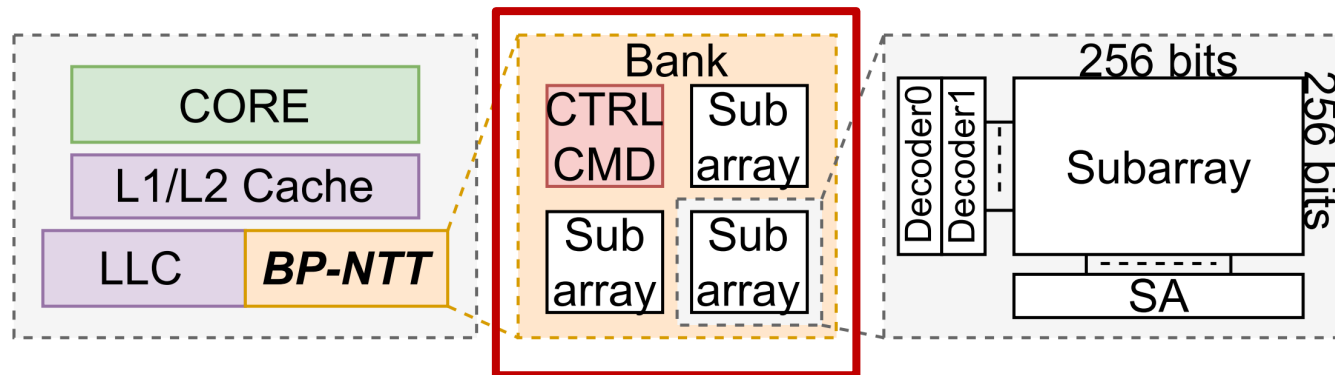
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*



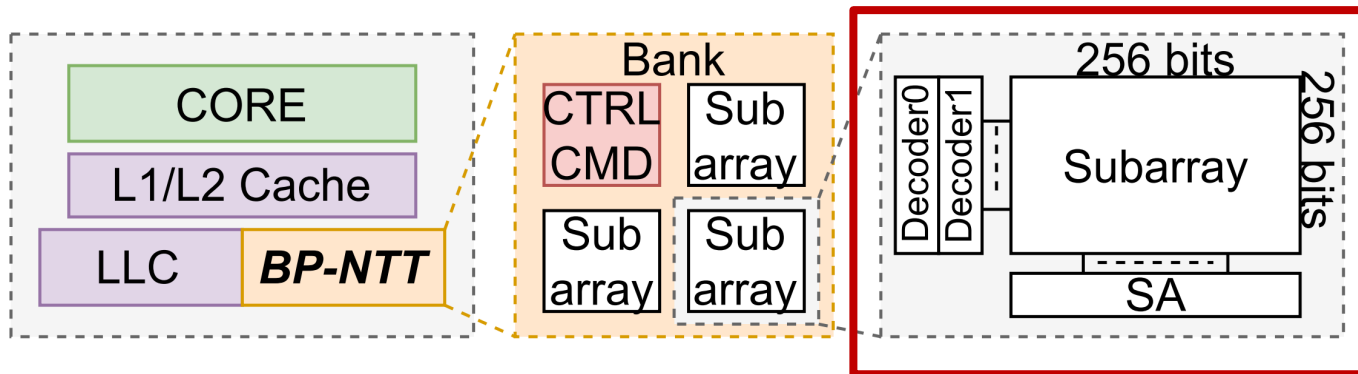
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*



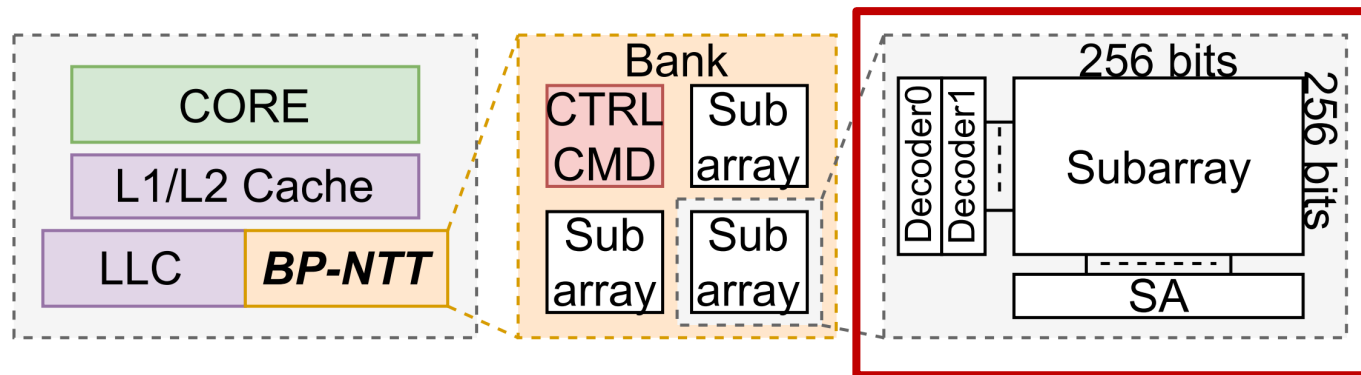
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*



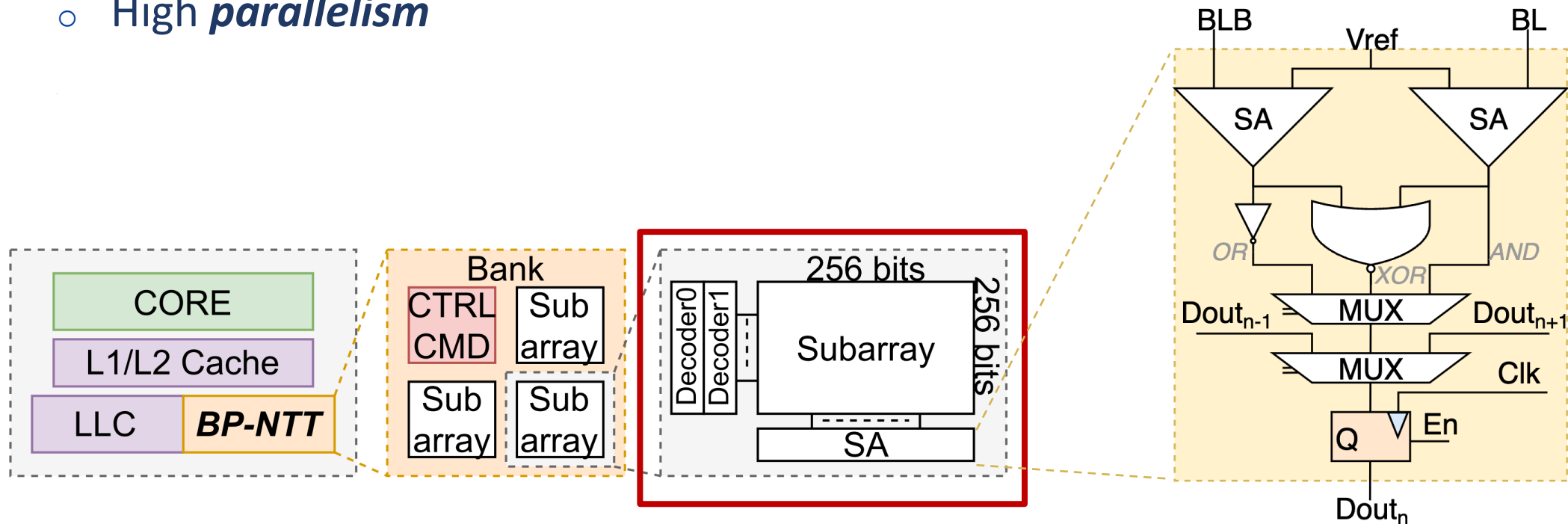
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*
  - High *parallelism*



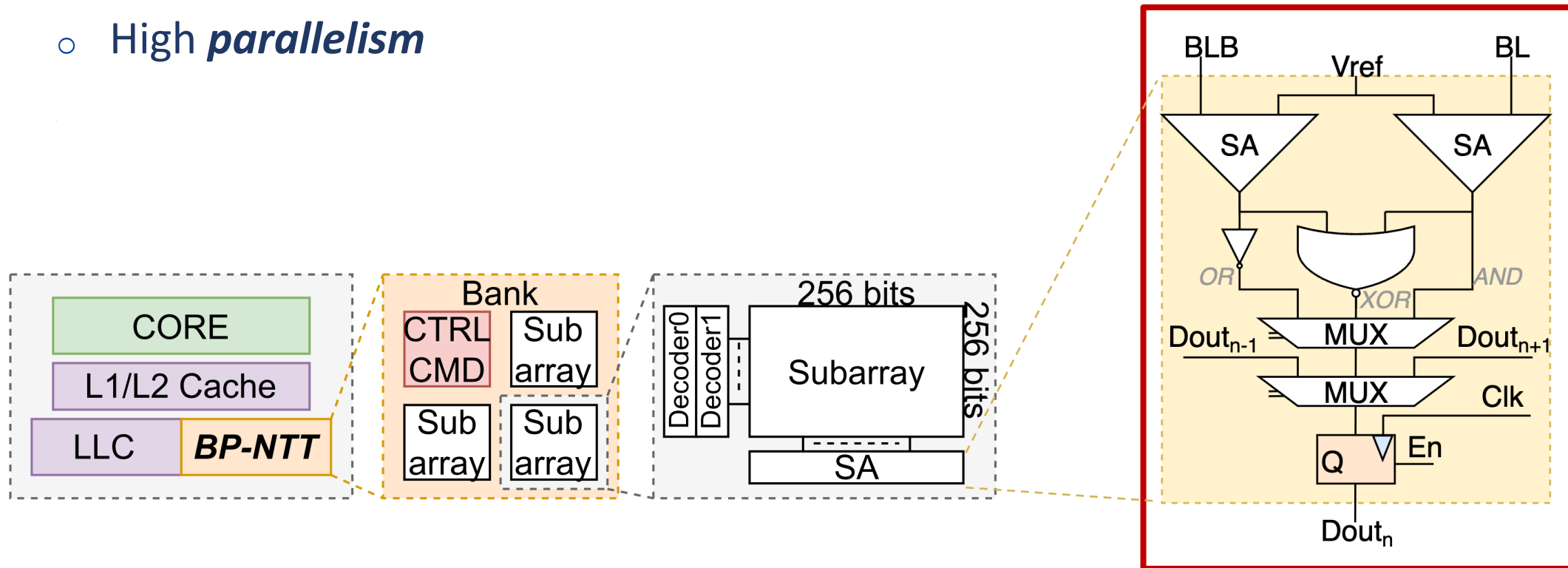
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*
  - High *parallelism*



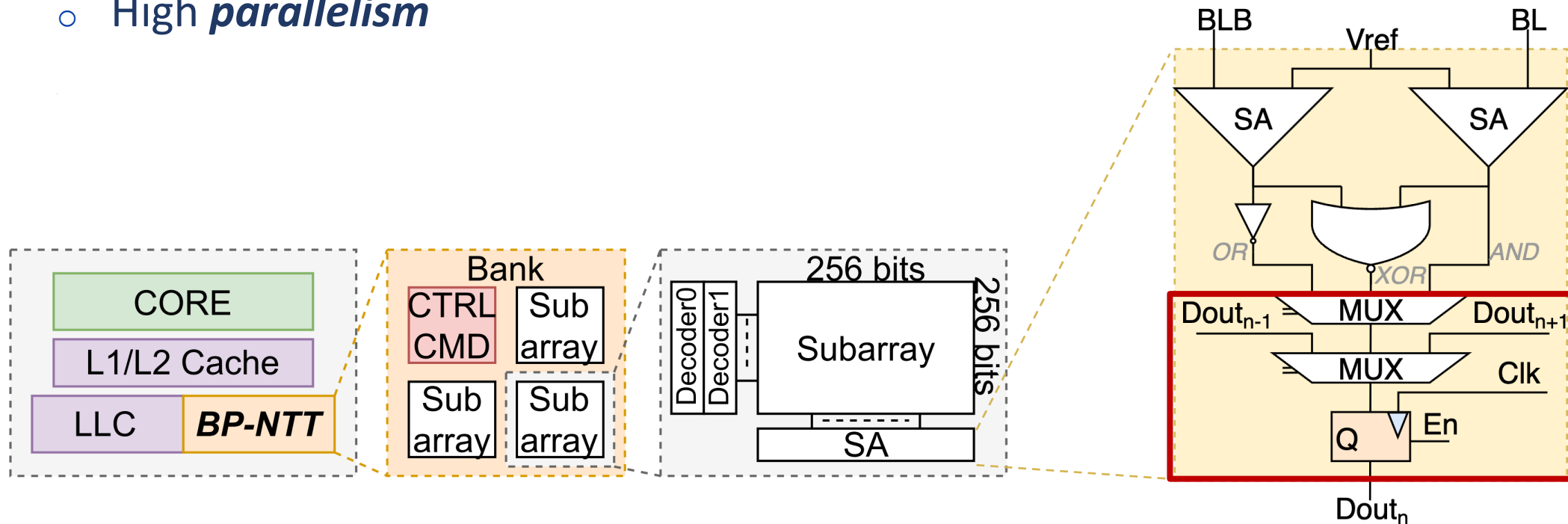
# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*
  - High *parallelism*



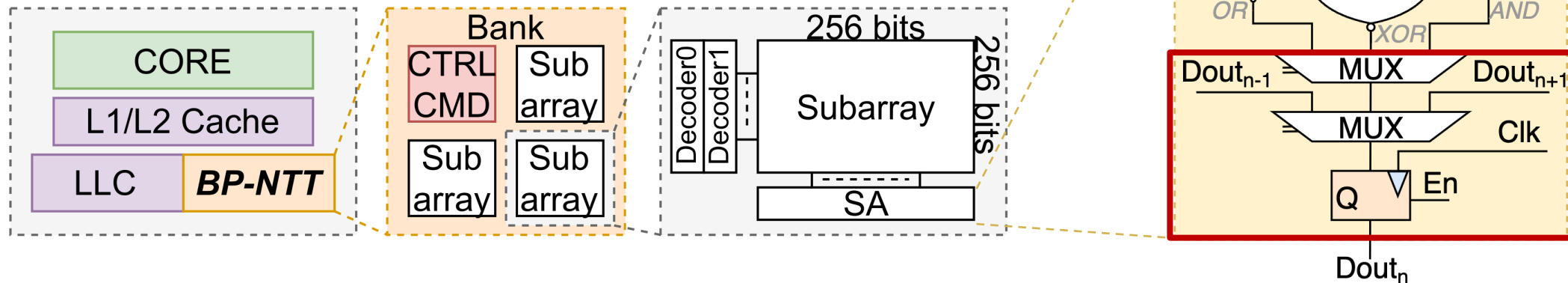
# BP-NTT: Repurposed LLC

- BP-NTT repurposes LLC to perform bitline computing [1]
  - High *security*
  - High *parallelism*



# BP-NTT: Repurposed LLC

- *BP-NTT* repurposes LLC to perform bitline computing [1]
  - High *security*
  - High *parallelism*
  - **< 2%** area overhead in 256x256 array





# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes LLC to perform bitline computing

High security 

High Throughput 

*BP-NTT* uses shift-optimized data alignment

Low Latency 

 Low Energy

*BP-NTT* employs bit-parallel modular multiplication

Low Latency 

 Small Area

# Motivation for Shift-optimization

---

# Motivation for Shift-optimization

---

- ~50% operations of 256-point 16-bit NTT are **shifting** in bit-serial

# Motivation for Shift-optimization

---

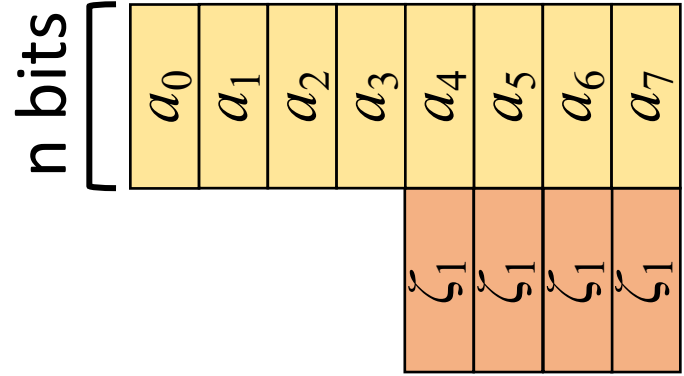
- ~50% operations of 256-point 16-bit NTT are **shifting** in bit-serial
- **Shifting** destroys parallelism due to **bit-by-bit** shift fashion

# ***BP-NTT***: Shift-optimized Data Alignment

---

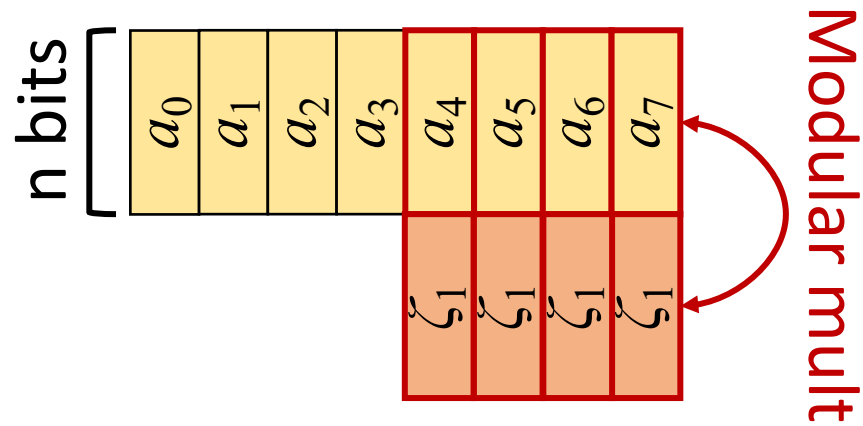
# BP-NTT: Shift-optimized Data Alignment

---



# BP-NTT: Shift-optimized Data Alignment

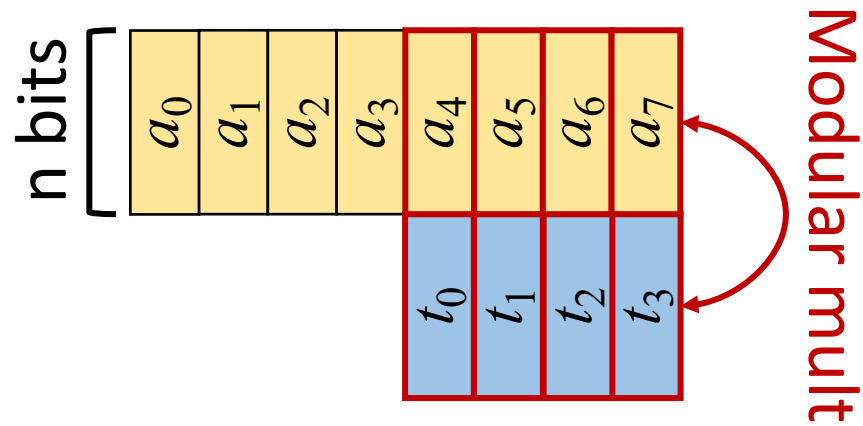
---



Step 1: Modular mult

# BP-NTT: Shift-optimized Data Alignment

---

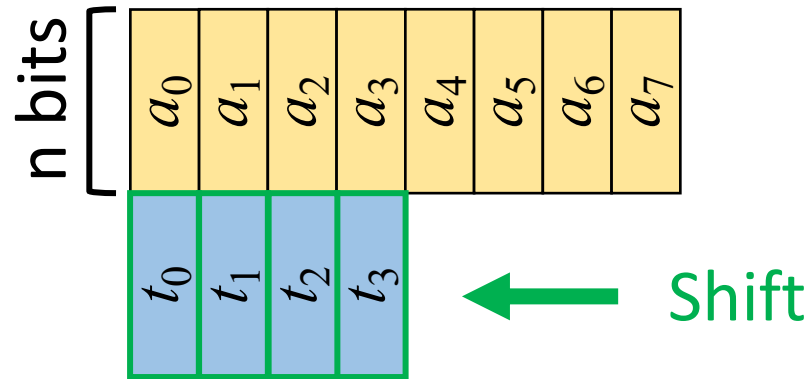


Step 1: Modular mult



# BP-NTT: Shift-optimized Data Alignment

---



Step 1: Modular mult

Step 2: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

---



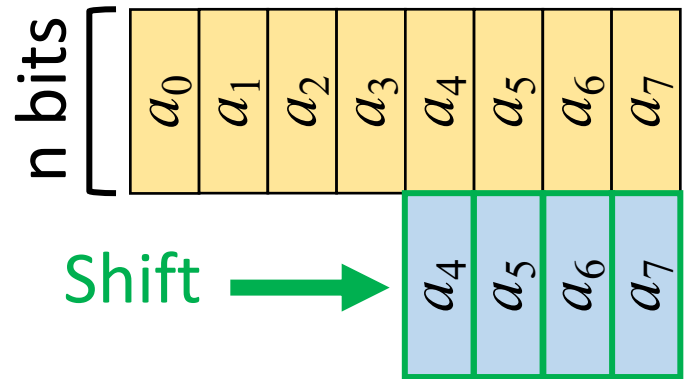
Step 1: Modular mult

Step 2: Shift & Write back

Step 3: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

---



Step 1: Modular mult

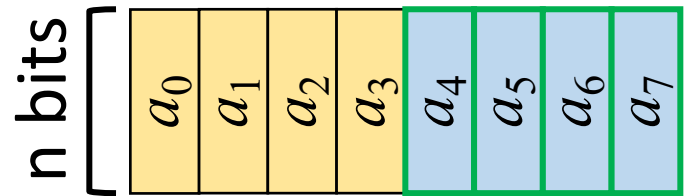
Step 2: Shift & Write back

Step 3: Add & Sub

Step 4: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

---



Write back

Step 1: Modular mult

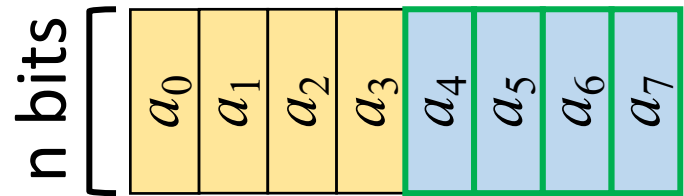
Step 2: Shift & Write back

Step 3: Add & Sub

Step 4: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Write back

Step 1: Modular mult



Step 2: Shift & Write back

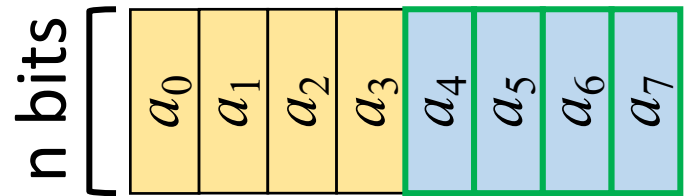
Step 3: Add & Sub



Step 4: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Write back



Step 1: Modular mult



Step 2: Shift & Write back

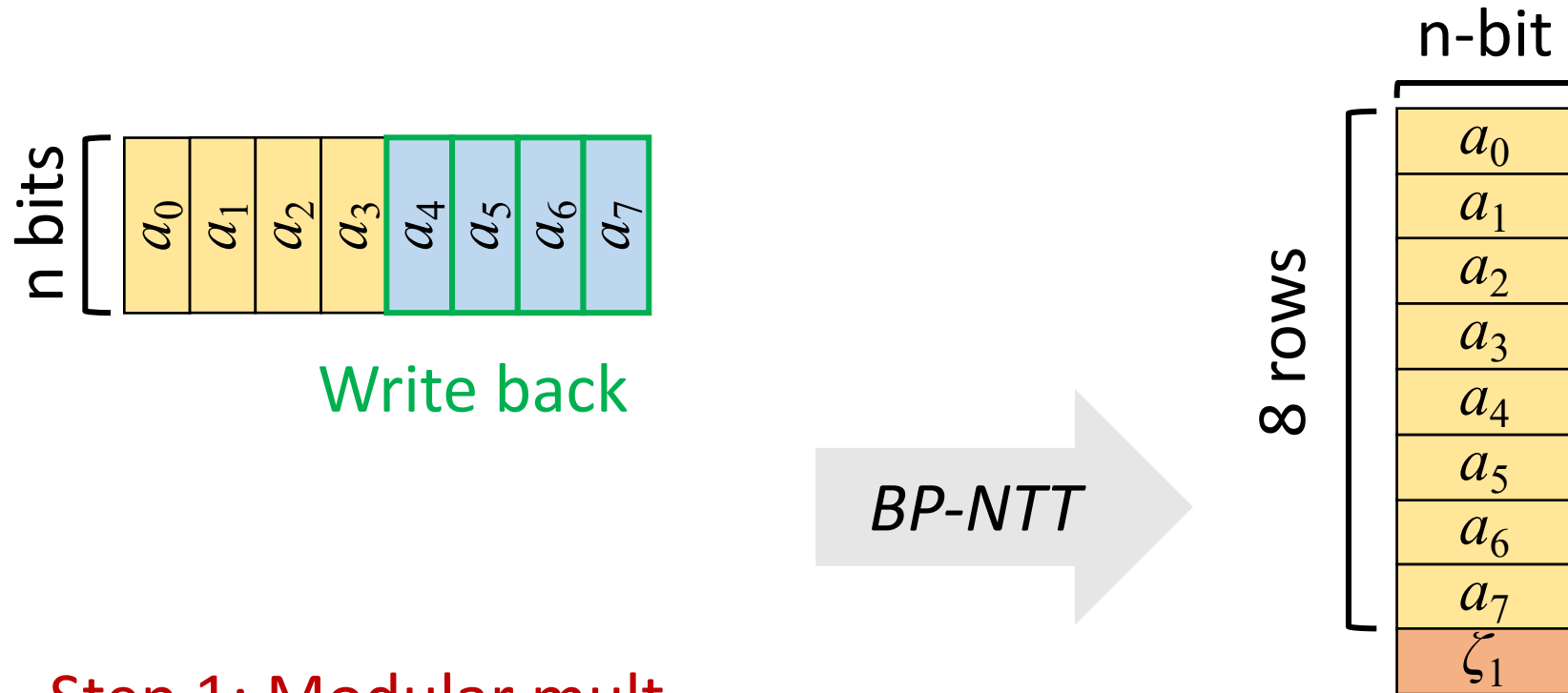
Step 3: Add & Sub



Step 4: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Step 1: Modular mult



Step 2: Shift & Write back

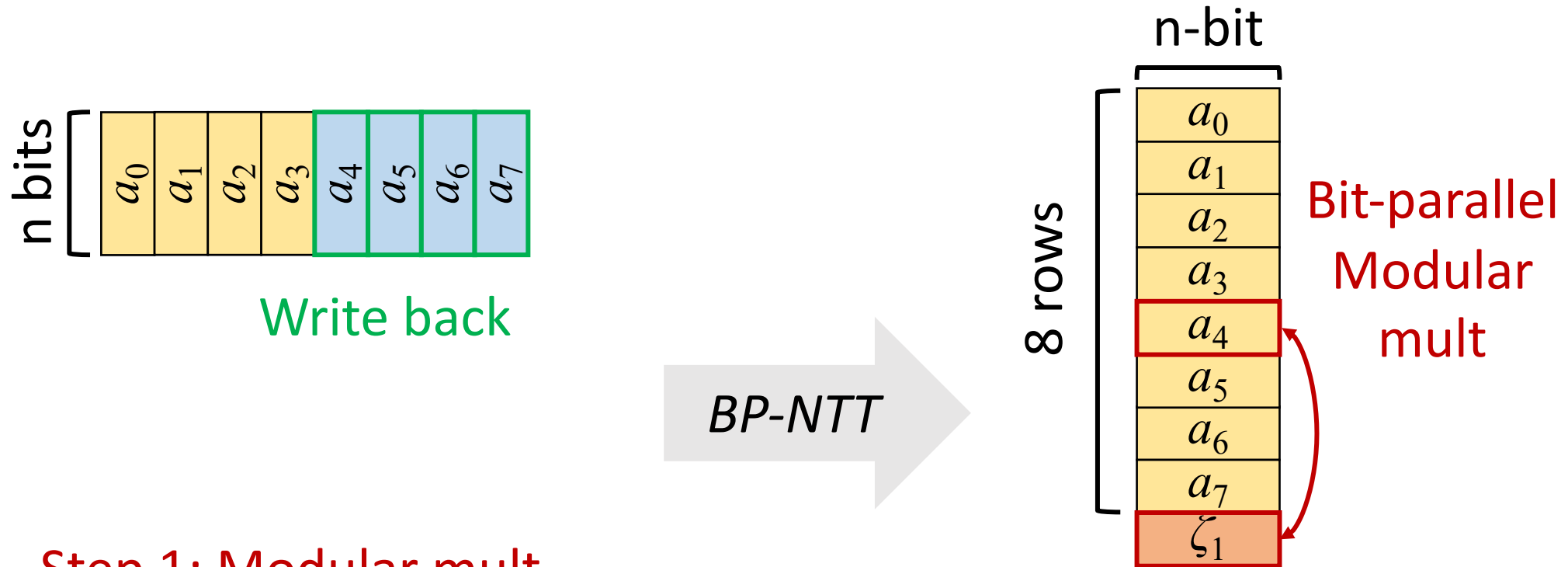
Step 3: Add & Sub



Step 4: Shift & Write back

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Step 1: Modular mult



Step 2: Shift & Write back

Step 3: Add & Sub



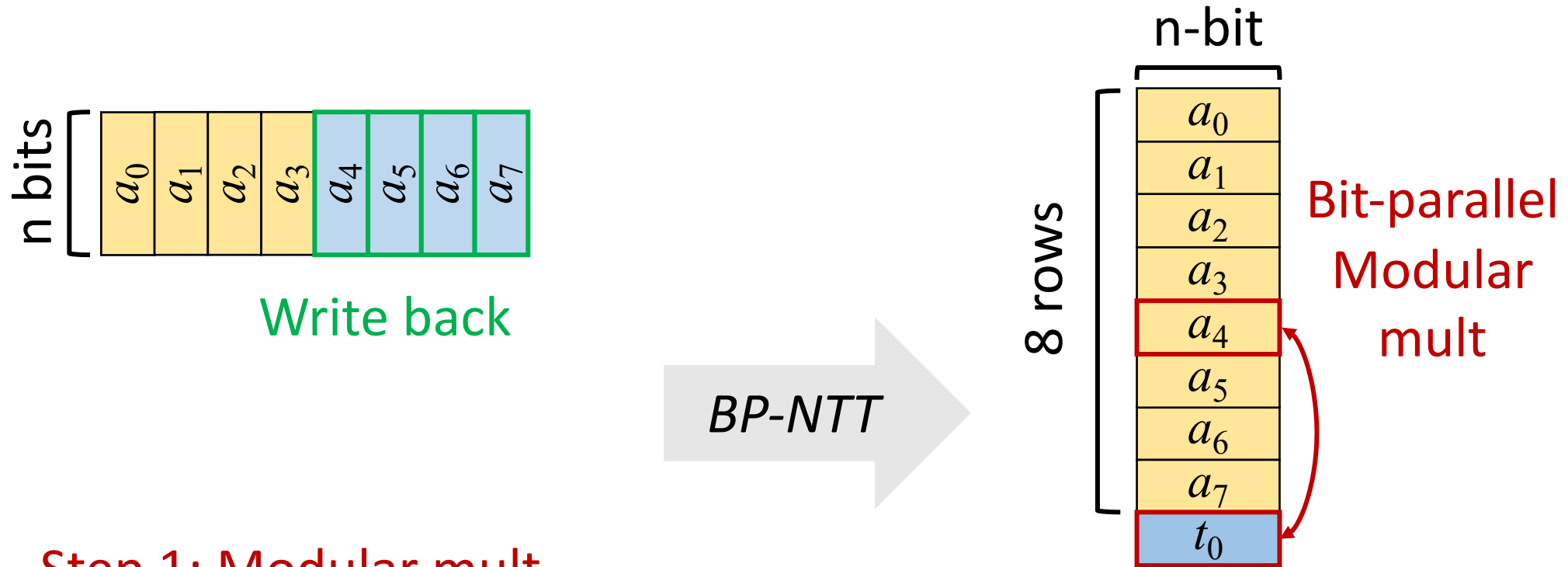
Step 4: Shift & Write back

Step 1: Modular mult



# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Step 1: Modular mult



Step 2: Shift & Write back

Step 3: Add & Sub

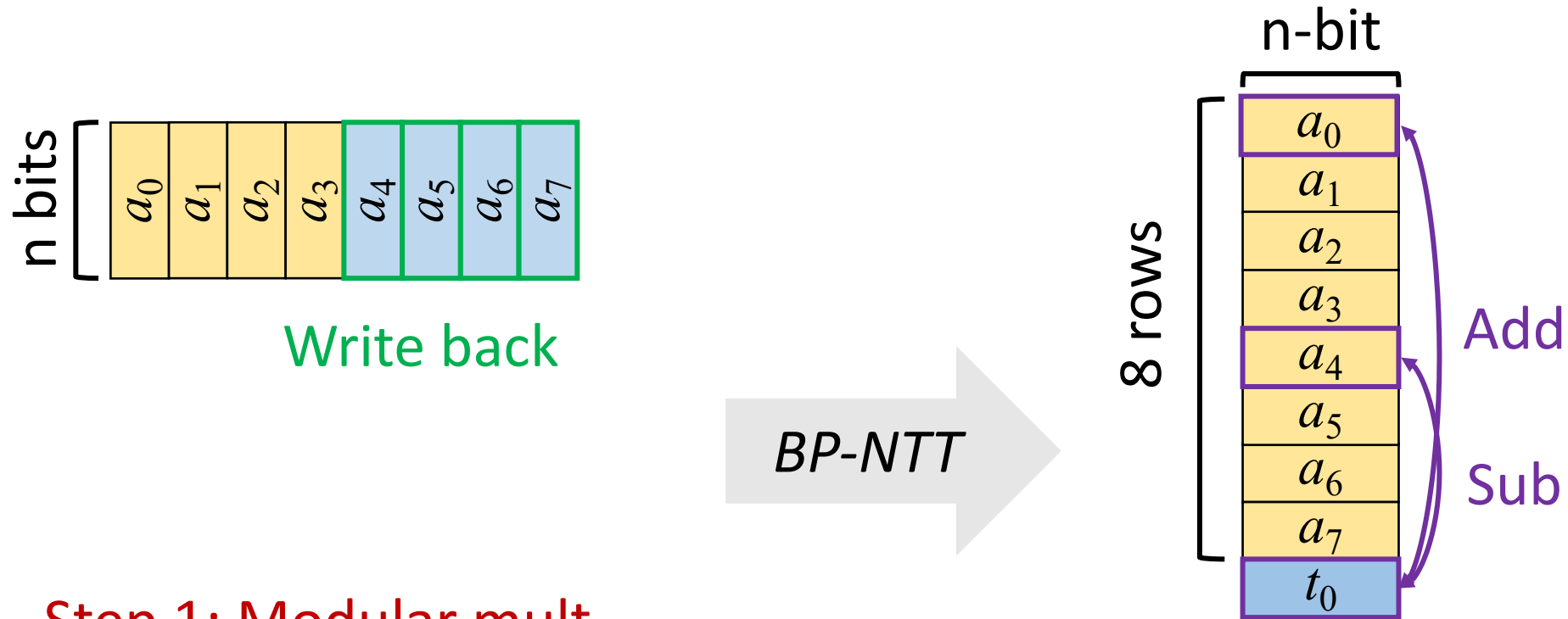


Step 4: Shift & Write back

Step 1: Modular mult

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Step 1: Modular mult



Step 2: Shift & Write back

Step 3: Add & Sub



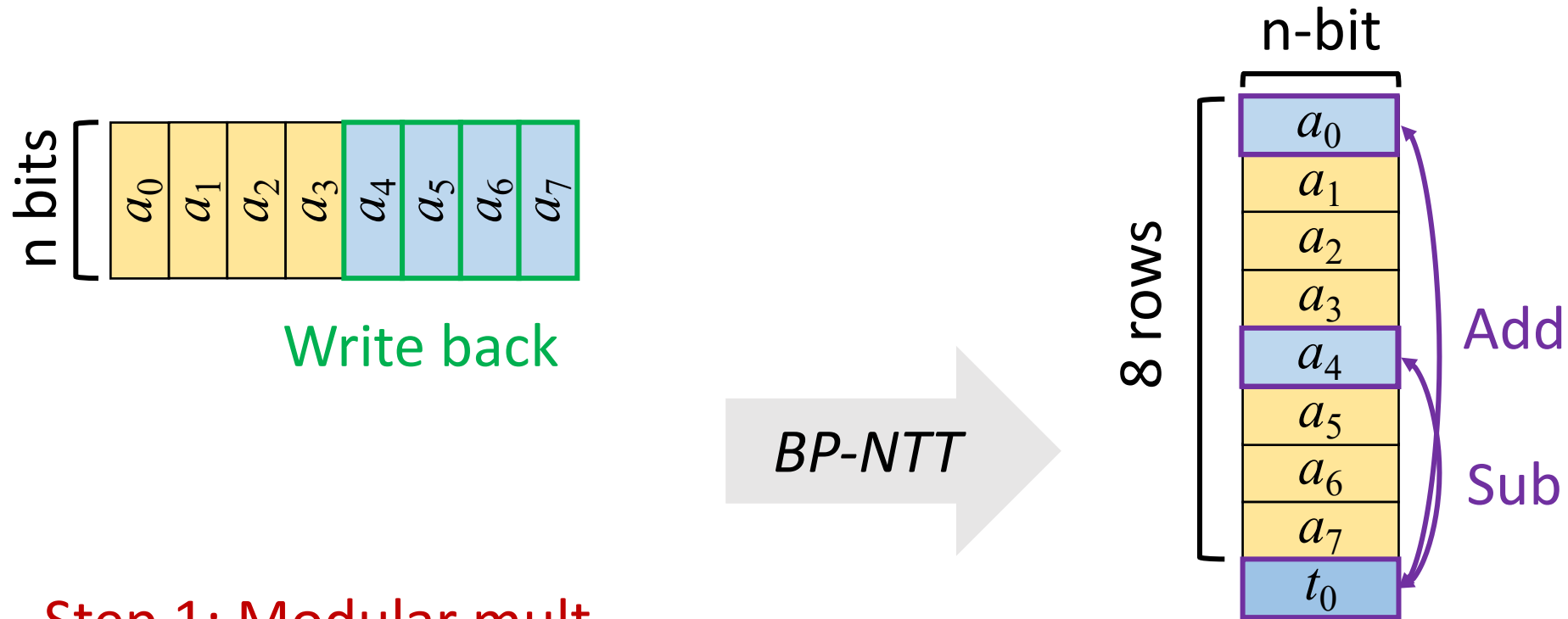
Step 4: Shift & Write back

Step 1: Modular mult

Step 2: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

- Inefficient data layout incurs *redundant shifts*



Step 1: Modular mult

Step 2: Shift & Write back

Step 3: Add & Sub

Step 4: Shift & Write back

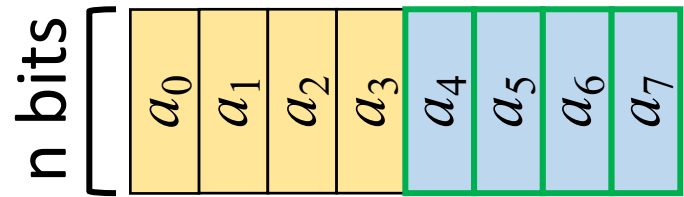


Step 1: Modular mult

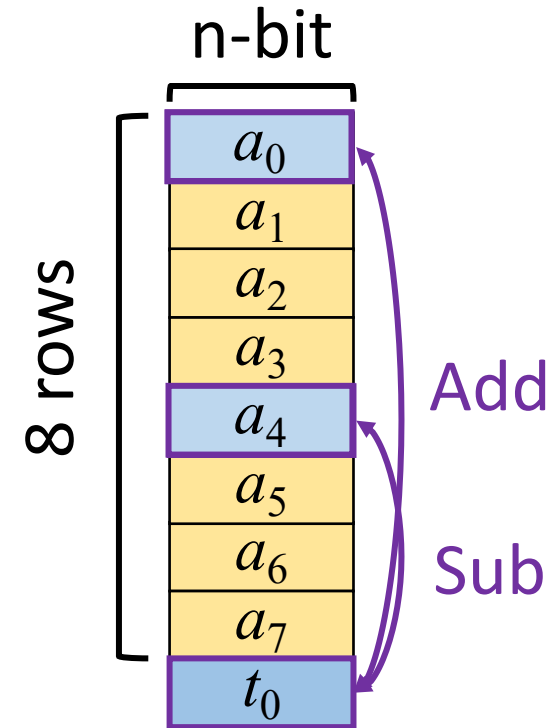
Step 2: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

- ❑ Inefficient data layout incurs *redundant shifts*
- ❑ **4-step stage** is simplified into **2-step**



Write back



Step 1: Modular mult

Step 2: Shift & Write back

Step 3: Add & Sub

Step 4: Shift & Write back

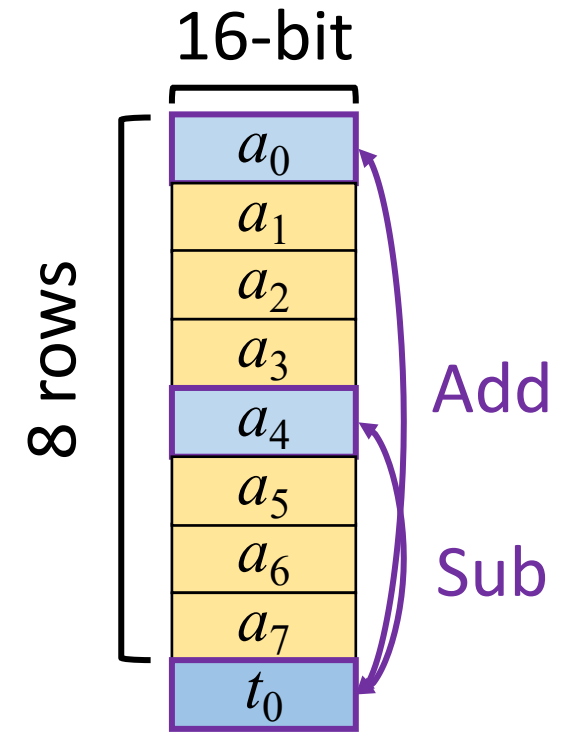


Step 1: Modular mult

Step 2: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

## □ Shift-optimized Data Alignment

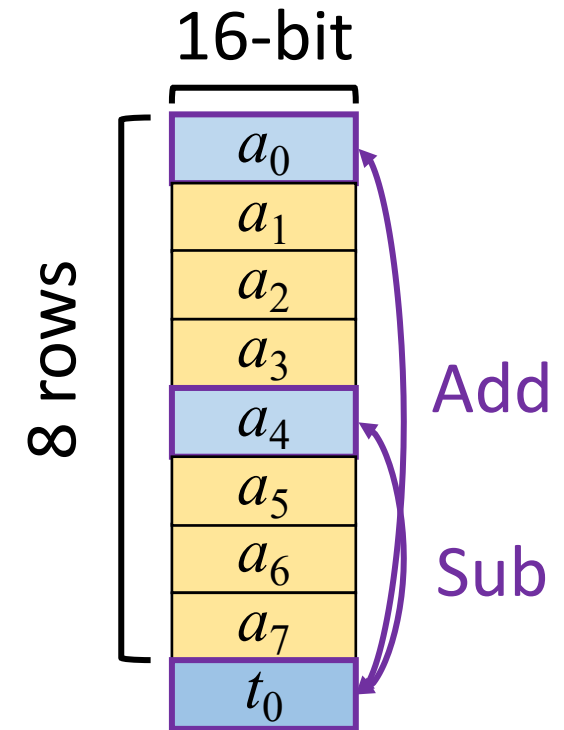


Step 1: Modular mult

Step 2: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

- Shift-optimized Data Alignment
  - Place coefficient per row



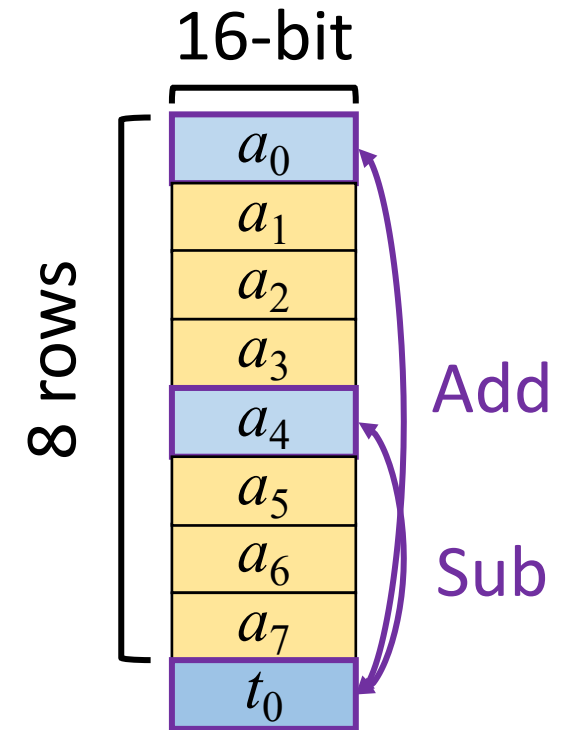
Step 1: Modular mult

Step 2: Add & Sub

# BP-NTT: Shift-optimized Data Alignment

## □ Shift-optimized Data Alignment

- **Place coefficient per row**
- **No shift operations** to align data



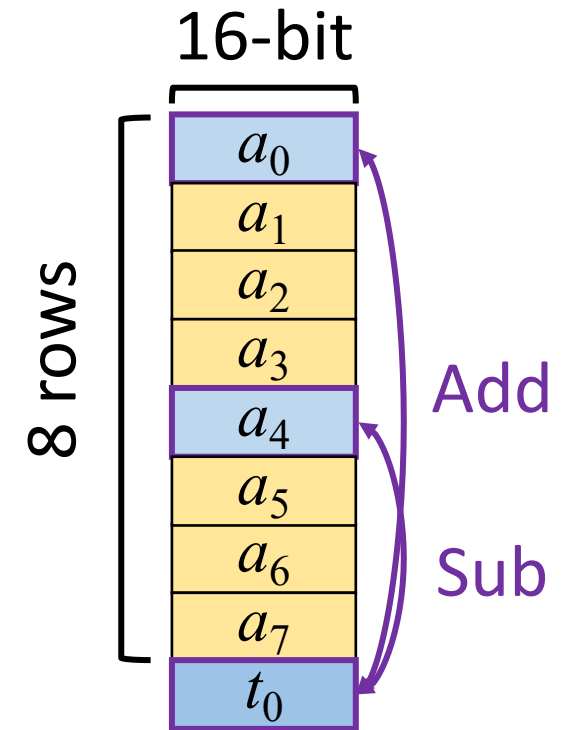
**Step 1: Modular mult**

**Step 2: Add & Sub**

# BP-NTT: Shift-optimized Data Alignment

## □ Shift-optimized Data Alignment

- **Place coefficient per row**
- **No shift operations** to align data
- Enable bit-parallel multiplication



**Step 1: Modular mult**

**Step 2: Add & Sub**



# Overview of Our Solution: *BP-NTT*

---

*BP-NTT* repurposes LLC to perform bitline computing

High security  

High Throughput  

*BP-NTT* uses shift-optimized data alignment

Low Latency  

 Low Energy 

*BP-NTT* employs bit-parallel modular multiplication

Low Latency  

 Small Area 

# Motivation for carry-save multiplication

---

# Motivation for carry-save multiplication

---

- Multiplication is based on multiple additions

# Motivation for carry-save multiplication

- Multiplication is based on multiple additions

$$\begin{array}{r} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \ 1 \ 1 \ 1 \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \times 1 \ 1 \ 1 \ 1 \\ \hline \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \ 1 \ 1 \ 1 \\ \phantom{+} \phantom{+} \phantom{+} + 1 \ 1 \ 1 \ 1 \\ \phantom{+} \phantom{+} + 1 \ 1 \ 1 \ 1 \\ + 1 \ 1 \ 1 \ 1 \\ \hline 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1 \end{array}$$

Multiplication

# Motivation for carry-save multiplication

- ❑ Multiplication is based on multiple additions
- ❑ Carry propagation ruins the parallelism from in-SRAM computing

$$\begin{array}{r} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \ 1 \ 1 \ 1 \\ \phantom{+} \phantom{+} \phantom{+} \times 1 \ 1 \ 1 \ 1 \\ \hline \phantom{+} \phantom{+} \phantom{+} 1 \ 1 \ 1 \ 1 \\ \phantom{+} + 1 \ 1 \ 1 \ 1 \\ + 1 \ 1 \ 1 \ 1 \\ \hline 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1 \end{array}$$

Multiplication

# Motivation for carry-save multiplication

---

- ❑ Multiplication is based on **multiple additions**
- ❑ **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity

$$\begin{array}{r} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ + \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ \hline 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1 \end{array}$$

Multiplication

# Motivation for carry-save multiplication

- Multiplication is based on **multiple additions**
- **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity

$$\begin{array}{r} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} \phantom{+} \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 1 \\ \hline 1 \phantom{+} 1 \phantom{+} 1 \phantom{+} 0 \phantom{+} 0 \phantom{+} 0 \phantom{+} 0 \phantom{+} 0 \phantom{+} 1 \end{array}$$

Multiplication

Complexity:  **$O(n^2)$**

# Motivation for carry-save multiplication

- ❑ Multiplication is based on **multiple additions**
- ❑ **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity

$$\begin{array}{r} \phantom{00} \phantom{00} 1\ 1\ 1\ 1 \\ \phantom{00} \phantom{00} \times 1\ 1\ 1\ 1 \\ \hline \phantom{00} \phantom{00} 1\ 1\ 1\ 1 \\ + \phantom{00} 1\ 1\ 1\ 1 \\ + \phantom{00} 1\ 1\ 1\ 1 \\ + \phantom{00} 1\ 1\ 1\ 1 \\ \hline 1\ 1\ 1\ 0\ 0\ 0\ 0\ 1 \end{array}$$

Multiplication  
Complexity:  $O(n^2)$

$$\begin{array}{r} \phantom{00} \phantom{00} 1\ 1\ 1\ 1 \\ \phantom{00} \phantom{00} + 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 1\ 1\ 1\ 0 \\ \phantom{00} \phantom{00} 0\ 0\ 0\ 1 \\ \text{Carry: } 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 0\ 1\ 1\ 0\ 0 \\ \phantom{00} \phantom{00} 0\ 0\ 0\ 1\ 0 \\ \text{Carry: } 0\ 0\ 0\ 1\ 0 \\ \hline \dots \end{array}$$

XOR  
AND  
XOR  
AND

n iters

Addition



# Motivation for carry-save multiplication

- ❑ Multiplication is based on **multiple additions**
- ❑ **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity

			1	1	1	1		
			×	1	1	1	1	
			1	1	1	1		
	+		1	1	1	1		
	+	1	1	1	1			
+	1	1	1	1				
1	1	1	0	0	0	0	1	

Multiplication  
Complexity:  **$O(n^2)$**

			1	1	1	1			
			+	0	0	0	1		
			Sum:	1	1	1	0	XOR	
				0	0	0	1	AND	
			Carry:	0	0	0	1		
			Sum:	0	1	1	0	XOR	
				0	0	0	1	AND	
			Carry:	0	0	0	1		
			.....						

} n iters

Addition Complexity:  **$O(n)$**

# Motivation for carry-save multiplication

- ❑ Multiplication is based on **multiple additions**
- ❑ **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity
  - Low parallelism (require extra columns for overflow bits)

				1	1	1	1				
				×	1	1	1	1			
				1	1	1	1				
			+	1	1	1	1				
		+	1	1	1	1					
	+	1	1	1	1						
				1	1	1	0	0	0	0	1

Multiplication  
Complexity: **O(n<sup>2</sup>)**

				1	1	1	1				
				+	0	0	0	1			
			Sum:	1	1	1	0		XOR		
				0	0	0	1		AND		
			Carry:	0	0	0	1				
			Sum:	0	1	1	0	0	XOR		
				0	0	0	1	0	AND		
			Carry:	0	0	0	1	0			
			.....								

} n iters

**Addition Complexity: O(n)**



# Motivation for carry-save multiplication

- ❑ Multiplication is based on **multiple additions**
- ❑ **Carry propagation** ruins the parallelism from in-SRAM computing
  - High computing complexity
  - Low parallelism (require extra columns for overflow bits)

				1	1	1	1	
			×	1	1	1	1	
			—	1	1	1	1	
		+		1	1	1	1	
	+			1	1	1	1	
+				1	1	1	1	
1	1	1	0	0	0	0	0	1

Multiplication  
Complexity:  **$O(n^2)$**

4-bit inputs  
 $\downarrow$  *mult*  
 8-bit result  
 $\downarrow$  *mod*  
 4-bit result

				1	1	1	1	
			+	0	0	0	1	
			—	1	1	1	0	XOR
			Sum:	1	1	1	0	AND
				0	0	0	1	
			Carry:	0	0	0	1	
			—	0	1	1	0	XOR
			Sum:	0	1	1	0	AND
				0	0	0	1	
			Carry:	0	0	0	1	
			—					
			.....					

Addition Complexity:  **$O(n)$**

} n iters

# ***BP-NTT:*** Bit-Parallel Modular Multiplication

---

# ***BP-NTT:*** Bit-Parallel Modular Multiplication

---

- Bit-parallel modular multiplication has *unnecessary carry propagations*

# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*

			1	1	1	1		
			+	0	0	0	1	
			<hr/>					
	Sum:		1	1	1	0	XOR	
			0	0	0	1	AND	
	Carry:		0	0	0	1		
			<hr/>					
	Sum:		0	1	1	0	0	XOR
			0	0	0	1	0	AND
	Carry:		0	0	0	1	0	
			<hr/>					
	.....						n iters	

Carry-propagation

Addition Complexity: **O(n)**

# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*

$$\begin{array}{r} \phantom{+} 1\ 1\ 1\ 1 \\ + 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 1\ 1\ 1\ 0 \text{ XOR} \\ \phantom{\text{Sum:}} 0\ 0\ 0\ 1 \text{ AND} \\ \text{Carry: } 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 0\ 1\ 1\ 0\ 0 \text{ XOR} \\ \phantom{\text{Sum:}} 0\ 0\ 0\ 1\ 0 \text{ AND} \\ \text{Carry: } 0\ 0\ 0\ 1\ 0 \\ \hline \end{array}$$

..... n iters

Carry-propagation

Addition Complexity:  **$O(n)$**





# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*

$$\begin{array}{r} \phantom{+} 1\ 1\ 1\ 1 \\ + 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 1\ 1\ 1\ 0 \text{ XOR} \\ \phantom{+} 0\ 0\ 0\ 1 \text{ AND} \\ \text{Carry: } 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 0\ 1\ 1\ 0\ 0 \text{ XOR} \\ \phantom{+} 0\ 0\ 0\ 1\ 0 \text{ AND} \\ \text{Carry: } 0\ 0\ 0\ 1\ 0 \\ \hline \end{array}$$

..... n iters

Carry-propagation

Addition Complexity:  **$O(n)$**

BP-NTT

$$\begin{array}{r} \phantom{+} 1\ 1\ 1\ 1 \\ + 0\ 0\ 0\ 1 \\ \hline \text{Sum: } 1\ 1\ 1\ 0 \text{ XOR} \\ \text{Carry: } 0\ 0\ 0\ 1 \text{ AND} \end{array}$$

Carry-save Addition

Complexity:  **$O(1)$**

# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*
- Multiplication  $O(n^2)$  is reduced into  $O(n)$  inspired by *carry-save addition*

$$\begin{array}{r}
 \phantom{+} 1\ 1\ 1\ 1 \\
 + 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 1\ 1\ 1\ 0 \quad \text{XOR} \\
 \phantom{\text{Sum:}} 0\ 0\ 0\ 1 \quad \text{AND} \\
 \text{Carry: } 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 0\ 1\ 1\ 0\ 0 \quad \text{XOR} \\
 \phantom{\text{Sum:}} 0\ 0\ 0\ 1\ 0 \quad \text{AND} \\
 \text{Carry: } 0\ 0\ 0\ 1\ 0 \\
 \hline
 \dots\dots \quad \text{n iters}
 \end{array}$$

Carry-propagation

Addition Complexity:  $O(n)$

BP-NTT

$$\begin{array}{r}
 \phantom{+} 1\ 1\ 1\ 1 \\
 + 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 1\ 1\ 1\ 0 \quad \text{XOR} \\
 \text{Carry: } 0\ 0\ 0\ 1 \quad \text{AND}
 \end{array}$$

Carry-save Addition

Complexity:  $O(1)$



# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*
- Multiplication  $O(n^2)$  is reduced into  $O(n)$  inspired by carry-save addition

$$\begin{array}{r}
 \phantom{+} 1\ 1\ 1\ 1 \\
 + 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 1\ 1\ 1\ 0 \quad \text{XOR} \\
 \phantom{+} 0\ 0\ 0\ 1 \quad \text{AND} \\
 \text{Carry: } 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 0\ 1\ 1\ 0\ 0 \quad \text{XOR} \\
 \phantom{+} 0\ 0\ 0\ 1\ 0 \quad \text{AND} \\
 \text{Carry: } 0\ 0\ 0\ 1\ 0 \\
 \hline
 \end{array}$$

.....

n iters

Carry-propagation

Addition Complexity:  $O(n)$



$$\begin{array}{r}
 \phantom{+} 1\ 1\ 1\ 1 \\
 + 0\ 0\ 0\ 1 \\
 \hline
 \text{Sum: } 1\ 1\ 1\ 0 \quad \text{XOR} \\
 \text{Carry: } 0\ 0\ 0\ 1 \quad \text{AND}
 \end{array}$$

4-bit inputs

4-bit results

Carry-save Addition

Complexity:  $O(1)$

# BP-NTT: Bit-Parallel Modular Multiplication

- Bit-parallel modular multiplication has *unnecessary carry propagations*
- Multiplication  $O(n^2)$  is reduced into  $O(n)$  inspired by carry-save addition

```

      1 1 1 1
    + 0 0 0 1
    -----
Sum:  1 1 1 0  XOR
      0 0 0 1  AND
Carry: 0 0 0 1
    -----
Sum:  0 1 1 0 0  XOR
      0 0 0 1 0  AND
Carry: 0 0 0 1 0
    -----
..... n iters

Carry-propagation
Addition Complexity:  $O(n)$ 

```



```

      1 1 1 1  4-bit inputs
    + 0 0 0 1
    -----
Sum:  1 1 1 0  XOR
Carry: 0 0 0 1  AND
    -----

      4-bit results

Carry-save Addition
Complexity:  $O(1)$ 

```

*Check paper for details*

# ***BP-NTT***: Overall Architecture

---

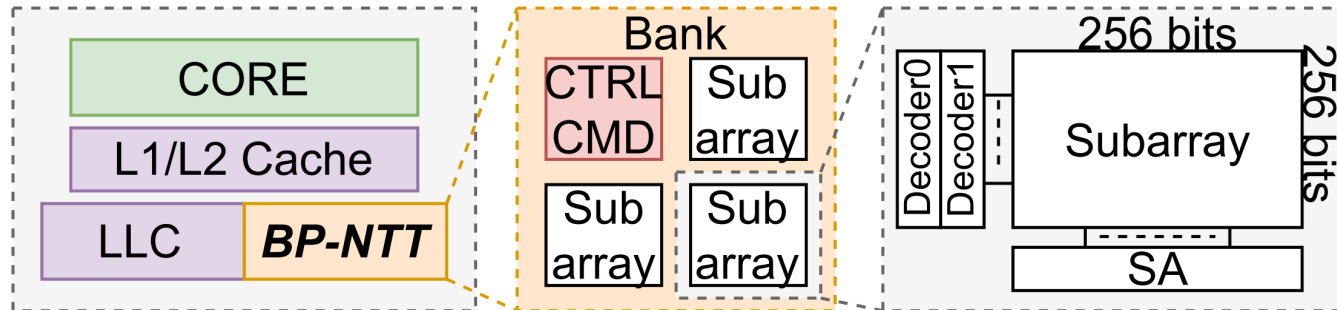
# ***BP-NTT***: Overall Architecture

---

High-performance, low-overhead, energy-efficient and flexible NTT engine

# BP-NTT: Overall Architecture

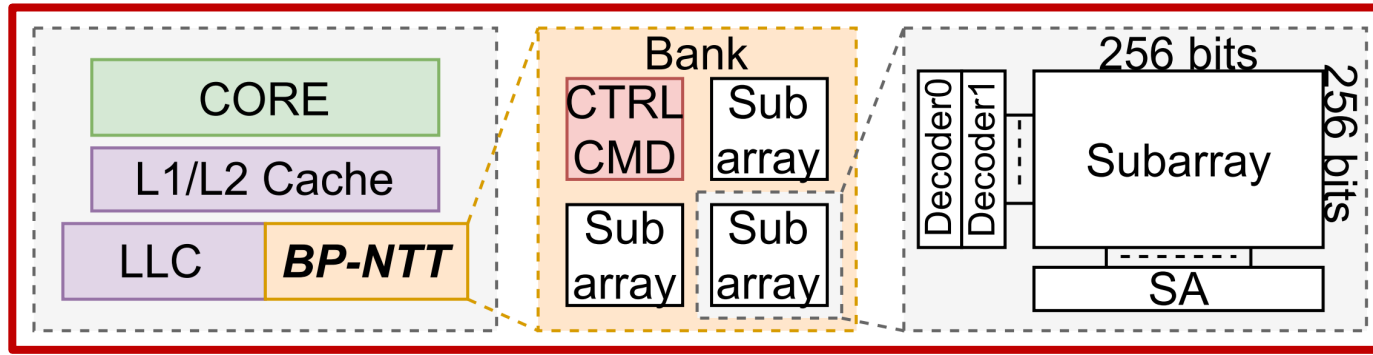
High-performance, low-overhead, energy-efficient and flexible NTT engine





# BP-NTT: Overall Architecture

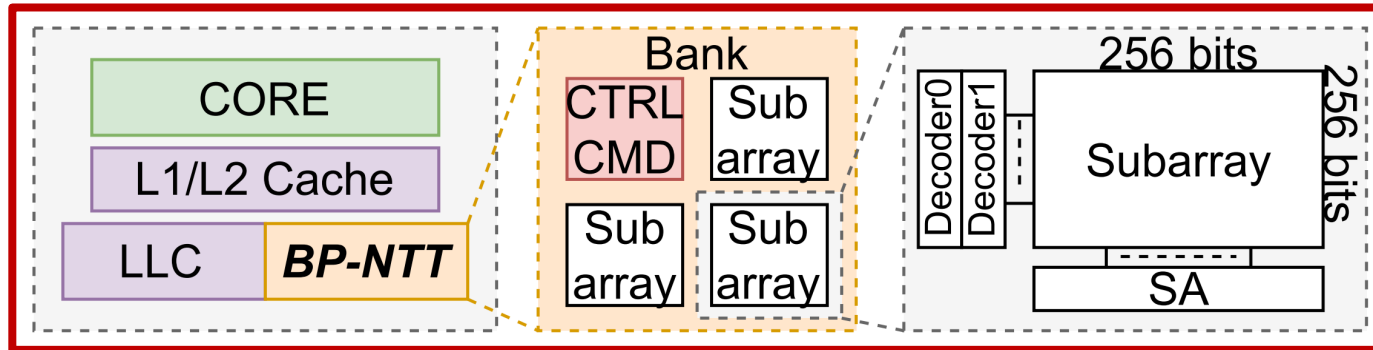
High-performance, low-overhead, energy-efficient and flexible NTT engine



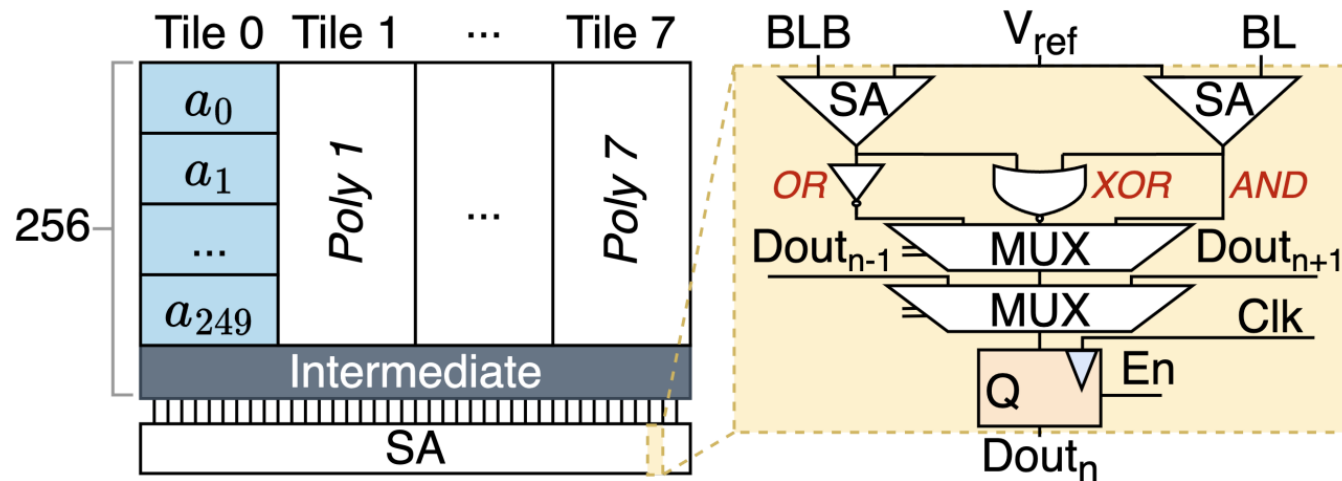
Security & Flexibility

# BP-NTT: Overall Architecture

High-performance, low-overhead, energy-efficient and flexible NTT engine

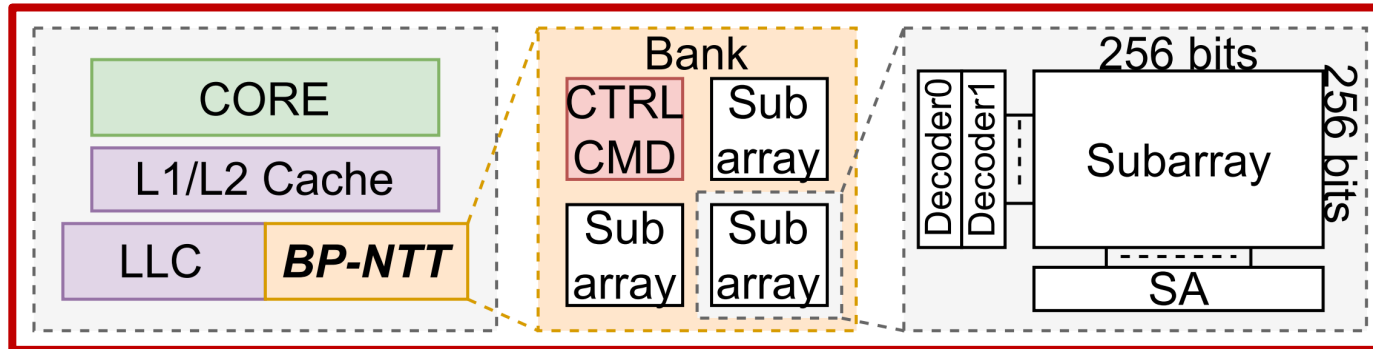


## Security & Flexibility

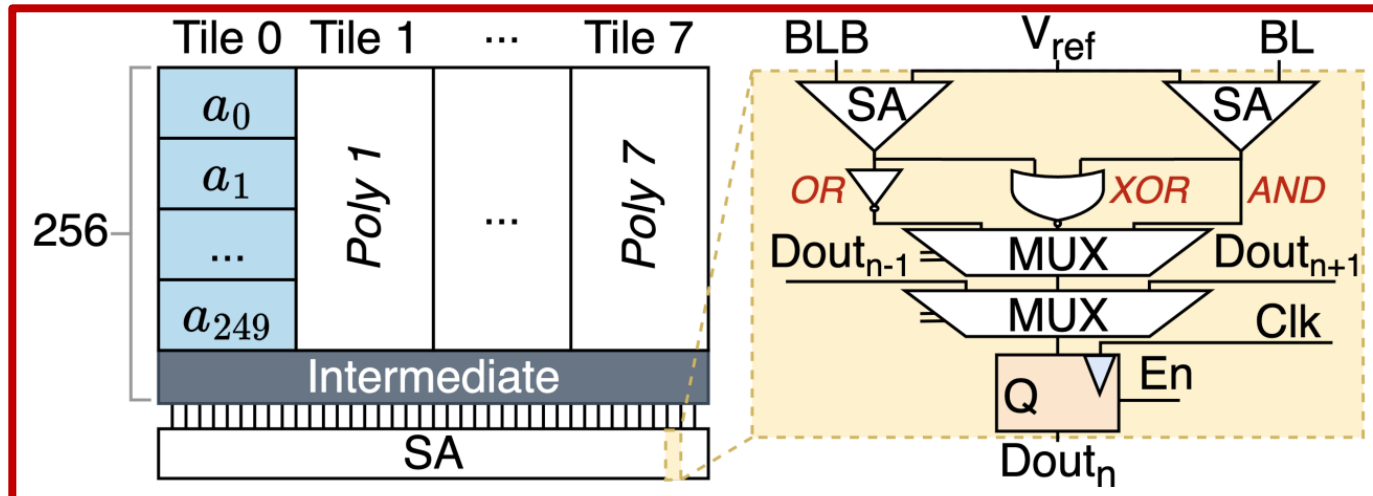


# BP-NTT: Overall Architecture

High-performance, low-overhead, energy-efficient and flexible NTT engine



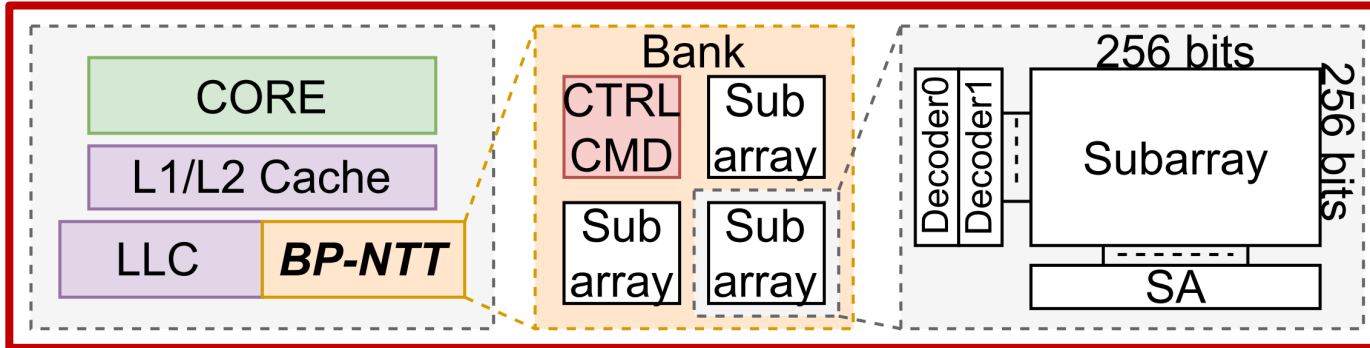
Security & Flexibility



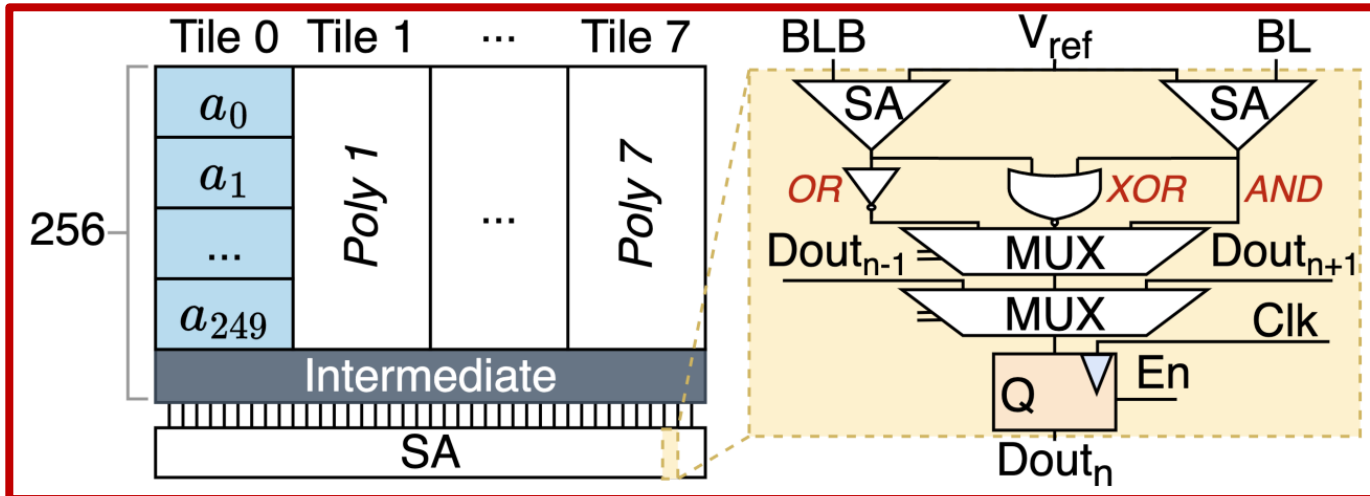
Throughput & Low-overhead

# BP-NTT: Overall Architecture

High-performance, low-overhead, energy-efficient and flexible NTT engine



Security & Flexibility



Throughput & Low-overhead

**Input:**  $n$ -bit  $A = (a_{n-1}, \dots, a_0)$ ,  $B = (b_{n-1}, \dots, b_0)$ ,  $M < R = 2^n$ , where  $n > 2$  and  $M \perp R$

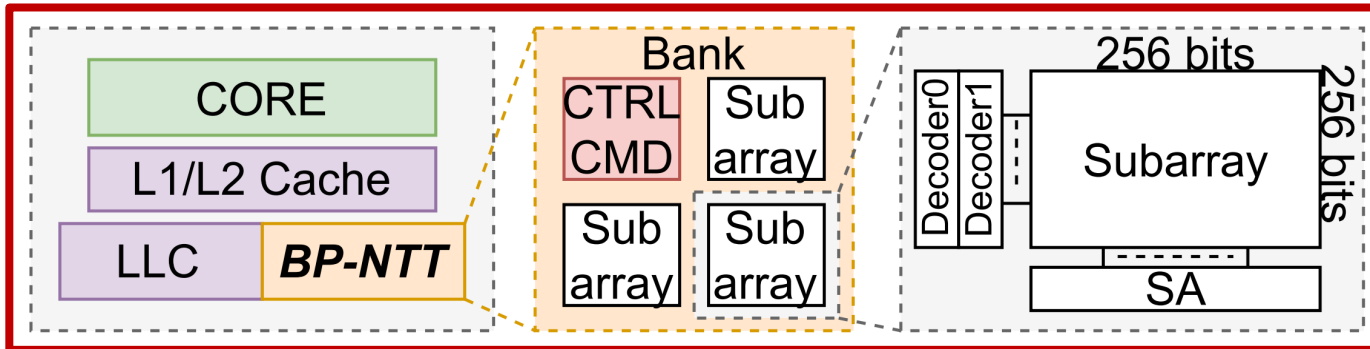
**Output:**  $ABR^{-1} \bmod M$

```

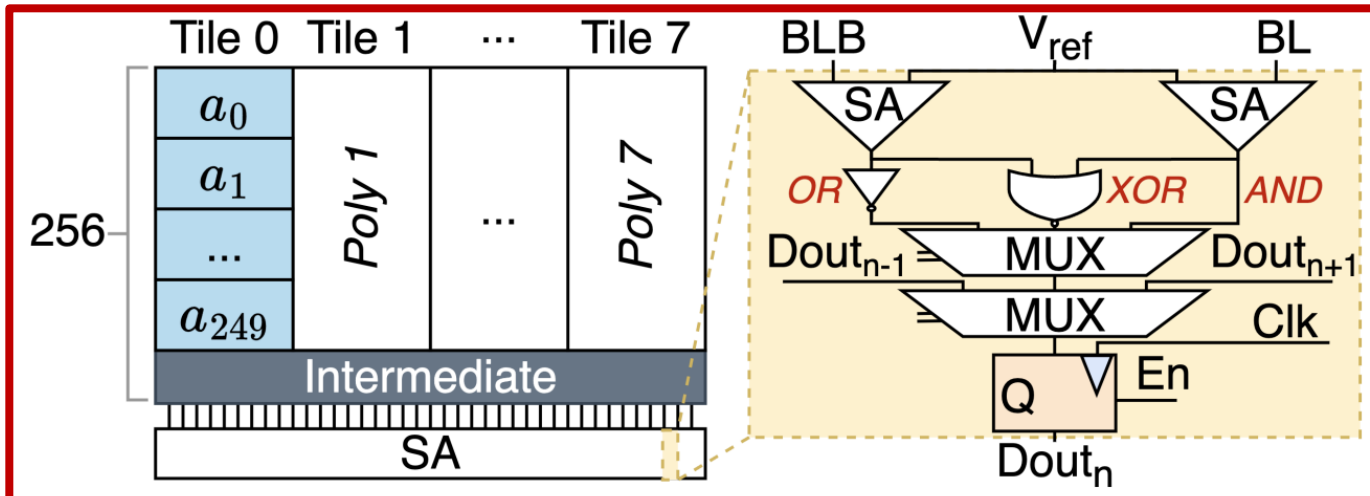
1: Sum := (s_{n-1}, ..., s_0) = 0 //Initialize
2: Carry := (c_{n-1}, ..., c_0) = 0
3: P := Sum + Carry << 1 // P = 0
4: for i = 0, n - 1 do
5:   if a_i == 1 then //Implicit compare
6:     c1, s1 = {Sum & B, Sum ⊕ B}
7:     Carry << 1 ← Observation 1
8:     c2, Sum = {Carry & s1, Carry ⊕ s1}
9:     Carry = c1 | c2 // P = P + a_i B
10:  end if
11:  m = (LSB(Sum) == 1) ? M : 0 // m = M or 0
12:  c1, s1 = {Sum & m, Sum ⊕ m}
13:  s1 >> 1 ← Observation 2
14:  c2, s2 = {s1 & c1, s1 ⊕ c1}
15:  c3, Sum = {Carry & s2, Carry ⊕ s2}
16:  Carry = c2 | c3 // P = P + m; P >> 1
17: end for
    
```

# BP-NTT: Overall Architecture

High-performance, low-overhead, energy-efficient and flexible NTT engine



## Security & Flexibility



## Throughput & Low-overhead

**Input:**  $n$ -bit  $A = (a_{n-1}, \dots, a_0)$ ,  $B = (b_{n-1}, \dots, b_0)$ ,  $M < R = 2^n$ , where  $n > 2$  and  $M \perp R$   
**Output:**  $ABR^{-1} \bmod M$

```

1: Sum :=  $(s_{n-1}, \dots, s_0) = 0$  // Initialize
2: Carry :=  $(c_{n-1}, \dots, c_0) = 0$ 
3:  $P := \text{Sum} + \text{Carry} \ll 1$  //  $P = 0$ 
4: for  $i = 0, n - 1$  do
5:   if  $a_i == 1$  then // Implicit compare
6:      $c1, s1 = \{\text{Sum} \& B, \text{Sum} \oplus B\}$ 
7:     Carry  $\ll 1$  // Observation 1
8:      $c2, \text{Sum} = \{\text{Carry} \& s1, \text{Carry} \oplus s1\}$ 
9:     Carry =  $c1 | c2$  //  $P = P + a_i B$ 
10:  end if
11:   $m = (\text{LSB}(\text{Sum}) == 1) ? M : 0$  //  $m = M$  or  $0$ 
12:   $c1, s1 = \{\text{Sum} \& m, \text{Sum} \oplus m\}$ 
13:   $s1 \gg 1$  // Observation 2
14:   $c2, s2 = \{s1 \& c1, s1 \oplus c1\}$ 
15:   $c3, \text{Sum} = \{\text{Carry} \& s2, \text{Carry} \oplus s2\}$ 
16:  Carry =  $c2 | c3$  //  $P = P + m; P \gg 1$ 
17: end for
    
```

## Latency

# Evaluation Methodology

---

- ❑ Tools:
  - PyMTL3 and OpenRAM for generating SRAM arrays
  - Synopsys Design Compiler for extracting latencies
  - Cadence Innovus for area consumption
- ❑ The array size of BP-NTT is  $256 \times 256$  following the ARM Cortex-M0+ microcontroller
- ❑ Area consumptions of in-ReRAM baselines are from DESTINY simulator
  - Optimistically estimated with only subarray area consumption excluding complex peripheral circuitry

# Comparison among Designs

---

# Comparison among Designs

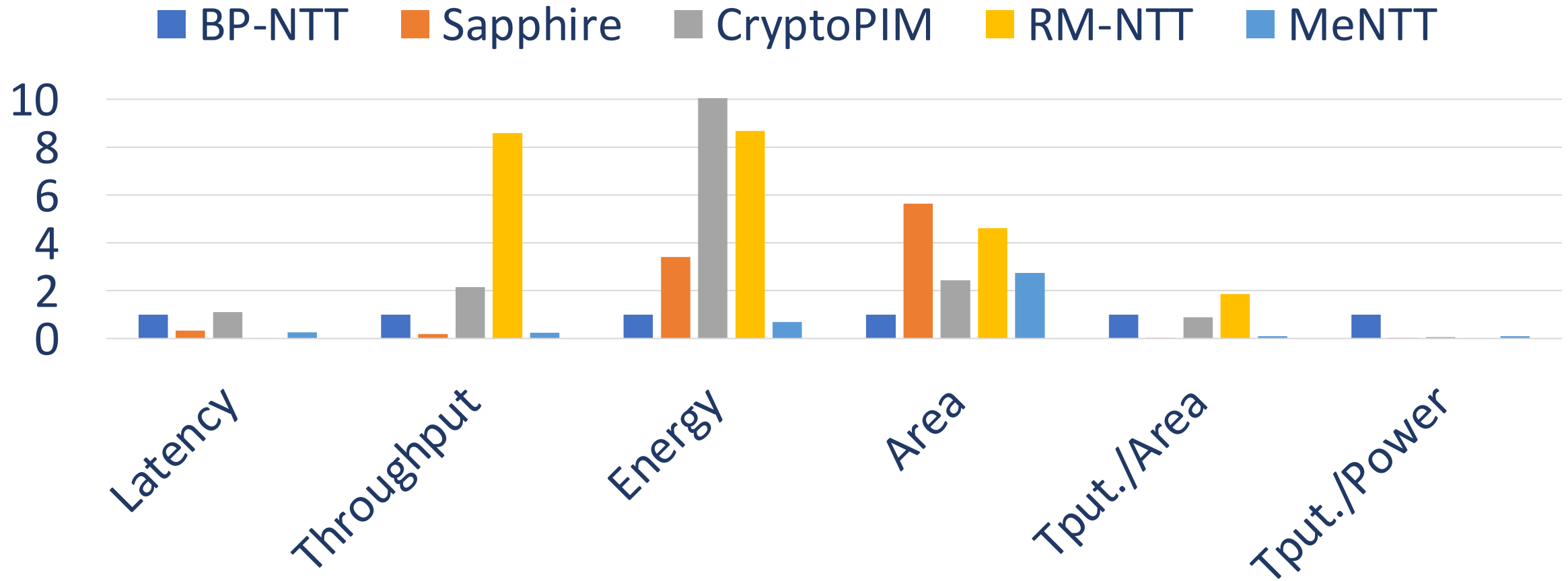
---

- Results are normalized to BP-NTT



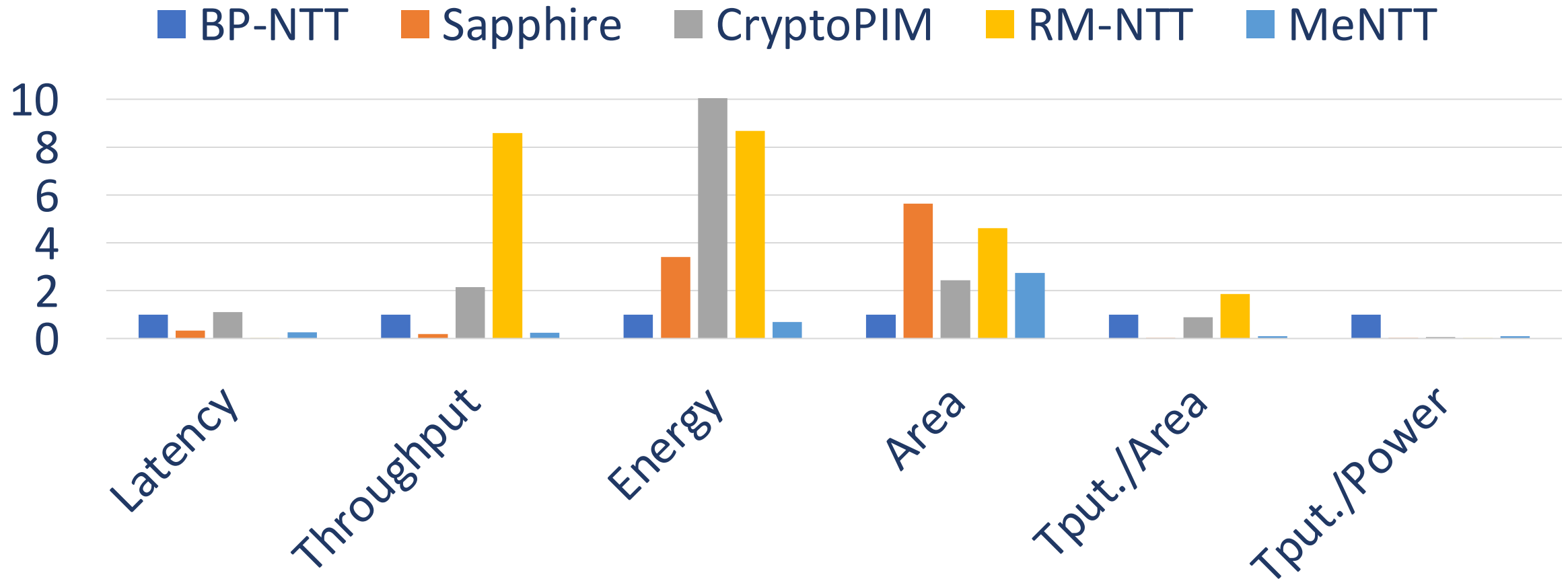
# Comparison among Designs

□ Results are normalized to BP-NTT



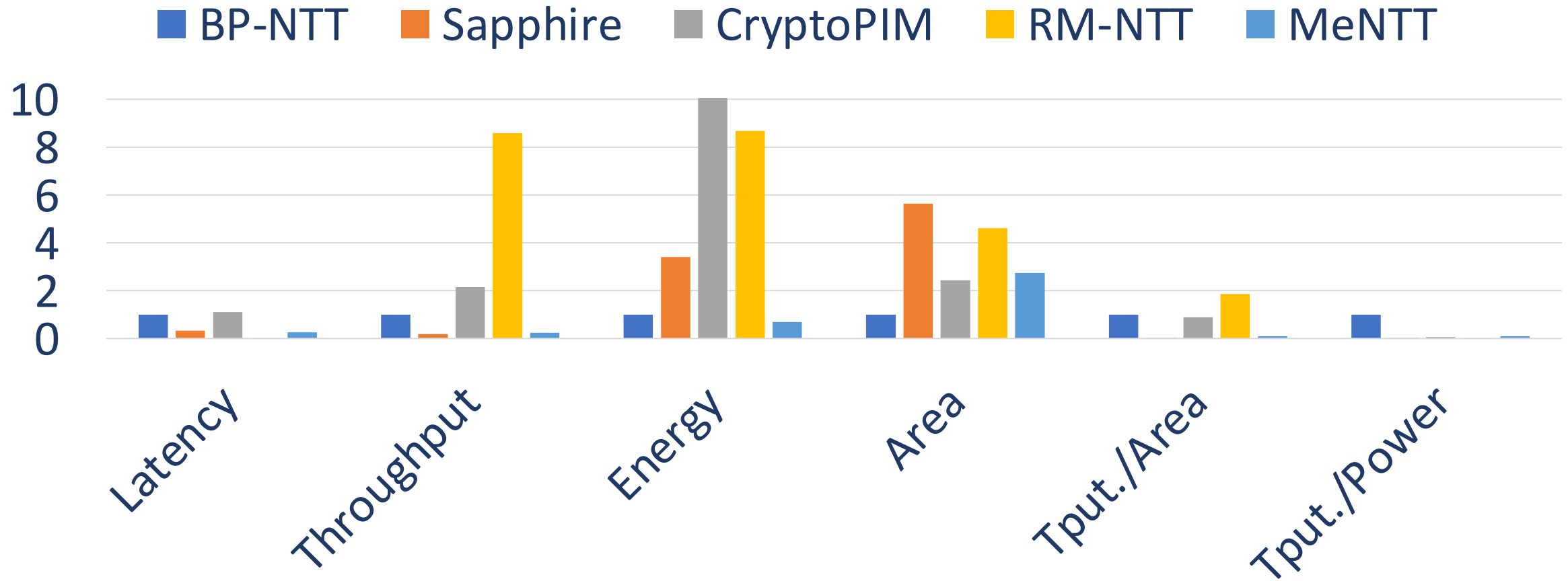
# Comparison among Designs

□ Results are normalized to BP-NTT



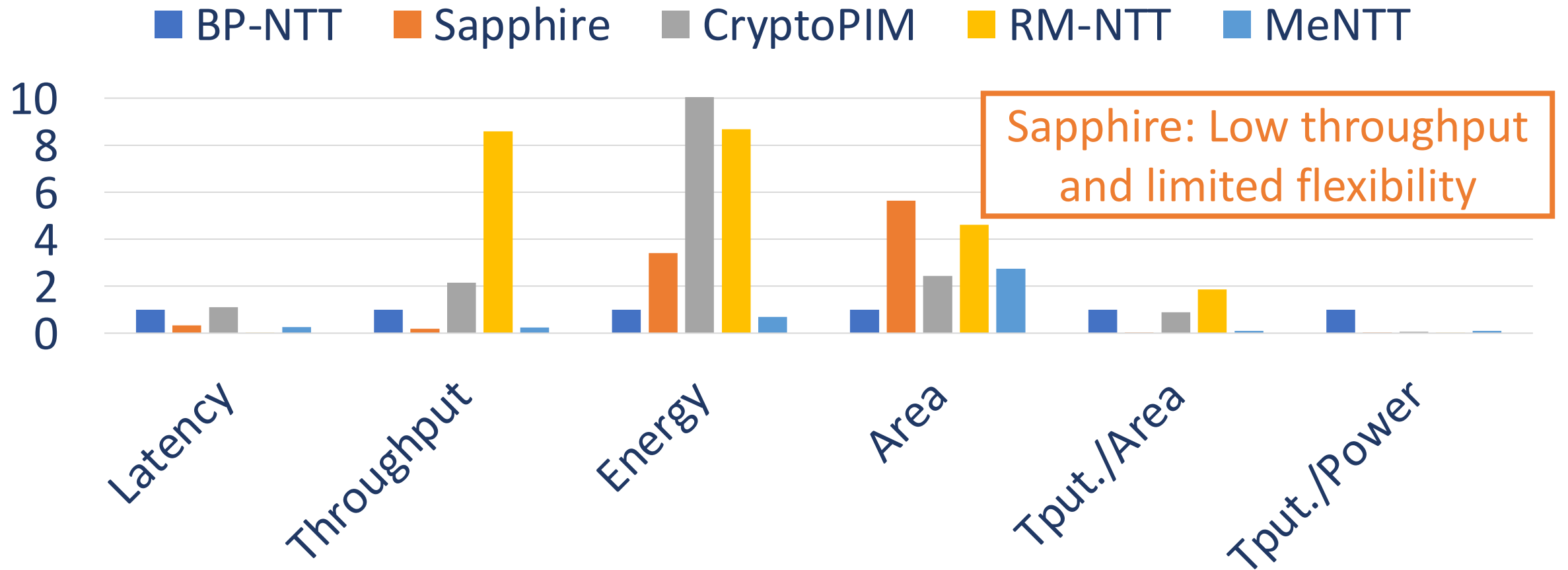
# Comparison among Designs

□ Results are normalized to BP-NTT



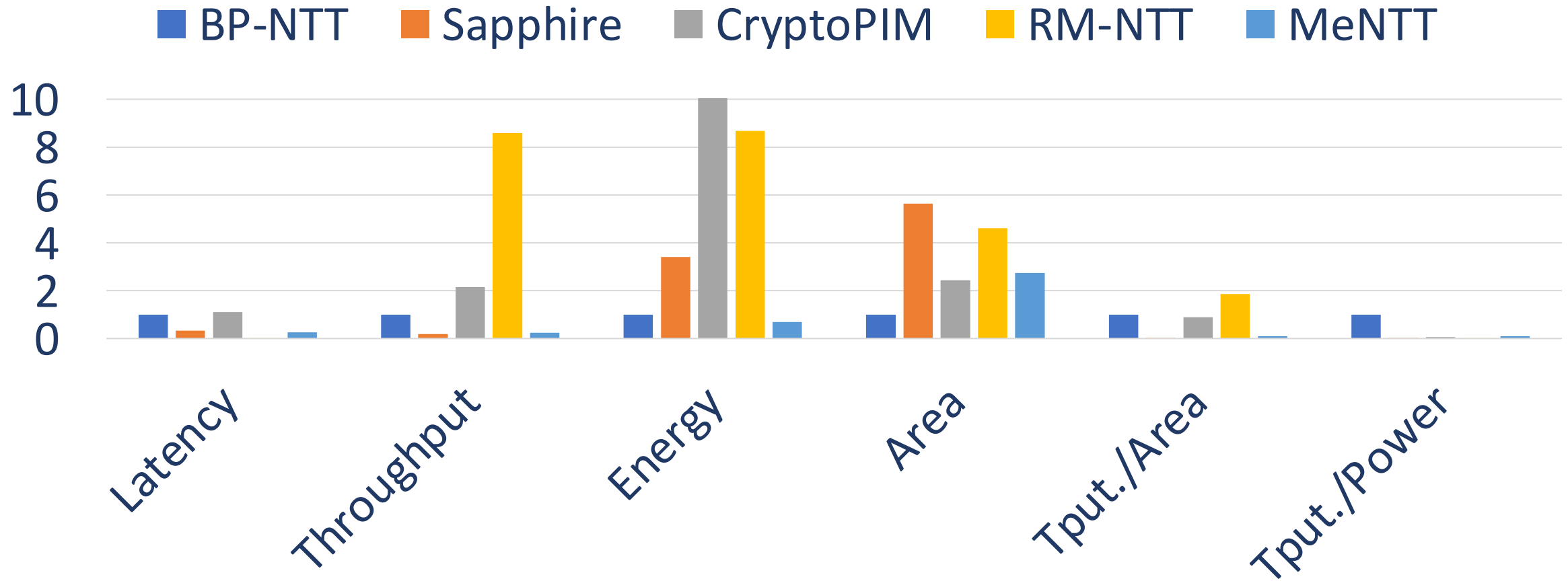
# Comparison among Designs

- Results are normalized to BP-NTT



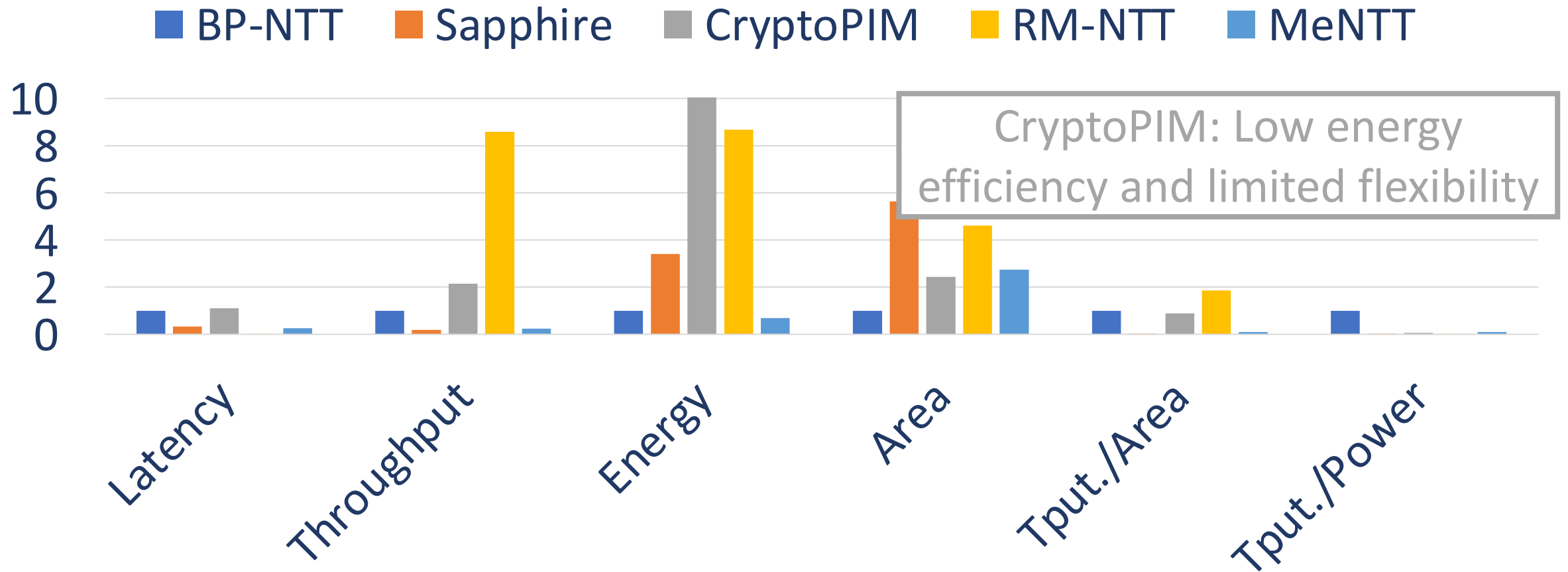
# Comparison among Designs

□ Results are normalized to BP-NTT



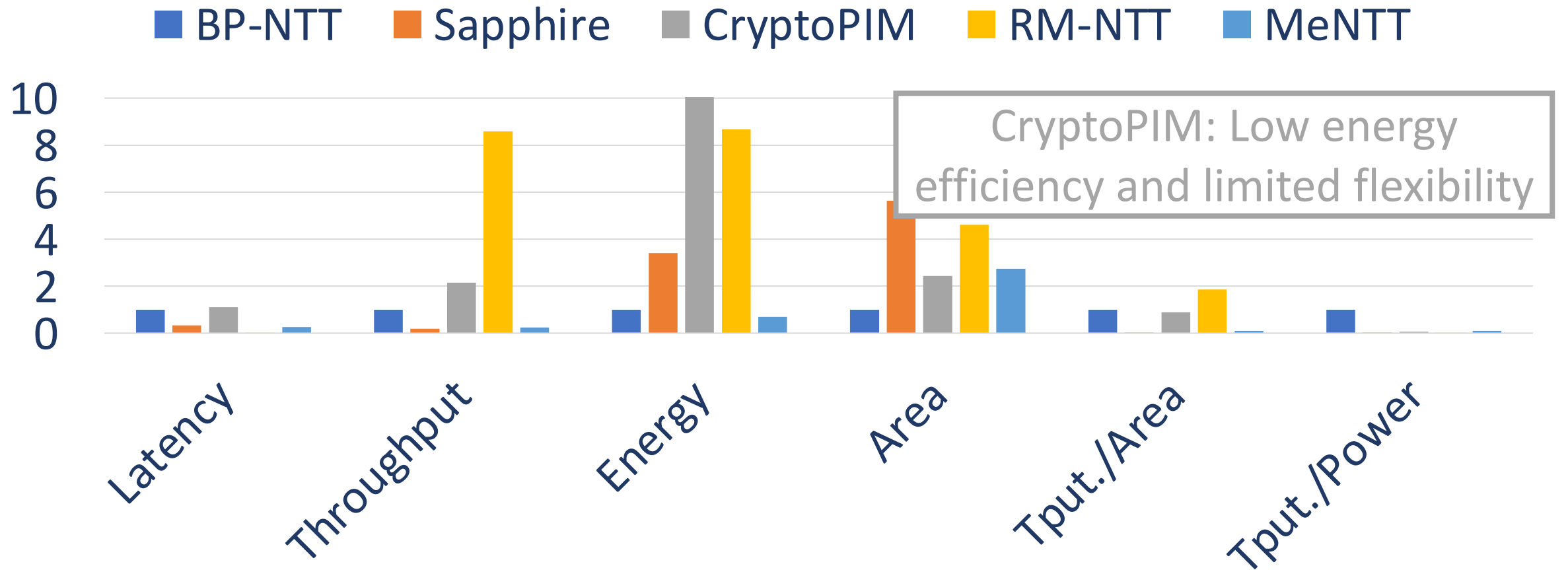
# Comparison among Designs

- Results are normalized to BP-NTT



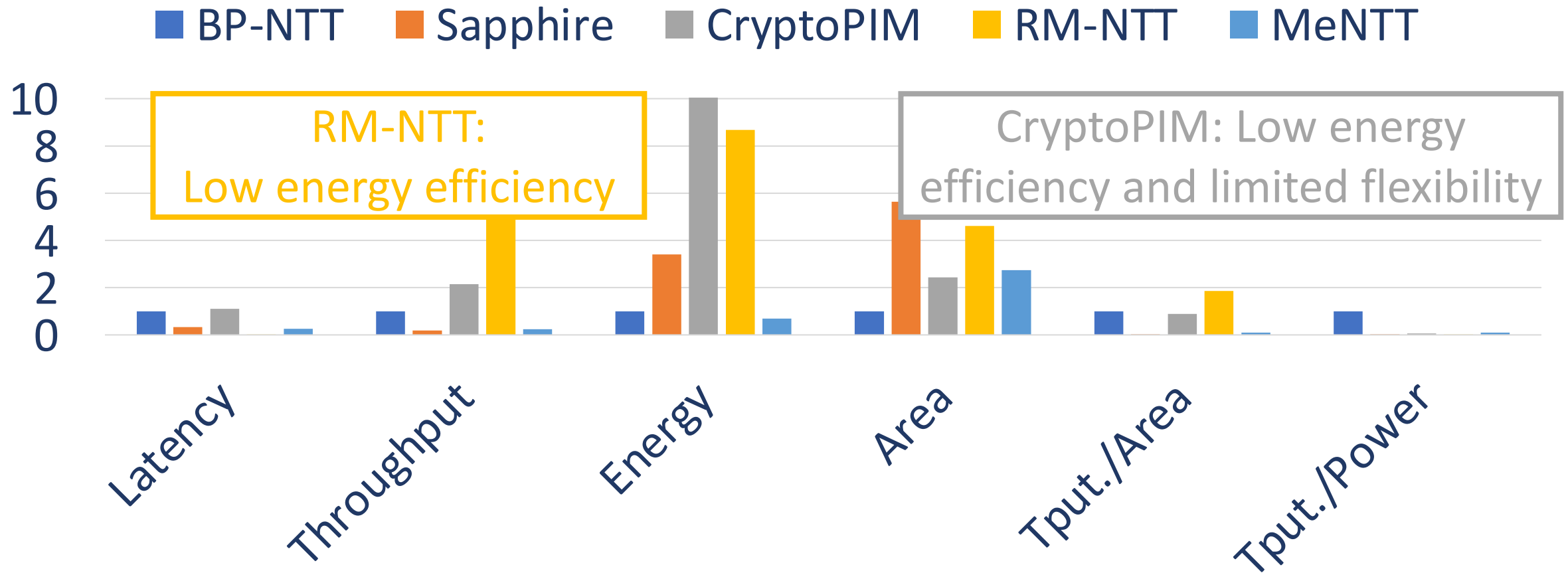
# Comparison among Designs

- Results are normalized to BP-NTT



# Comparison among Designs

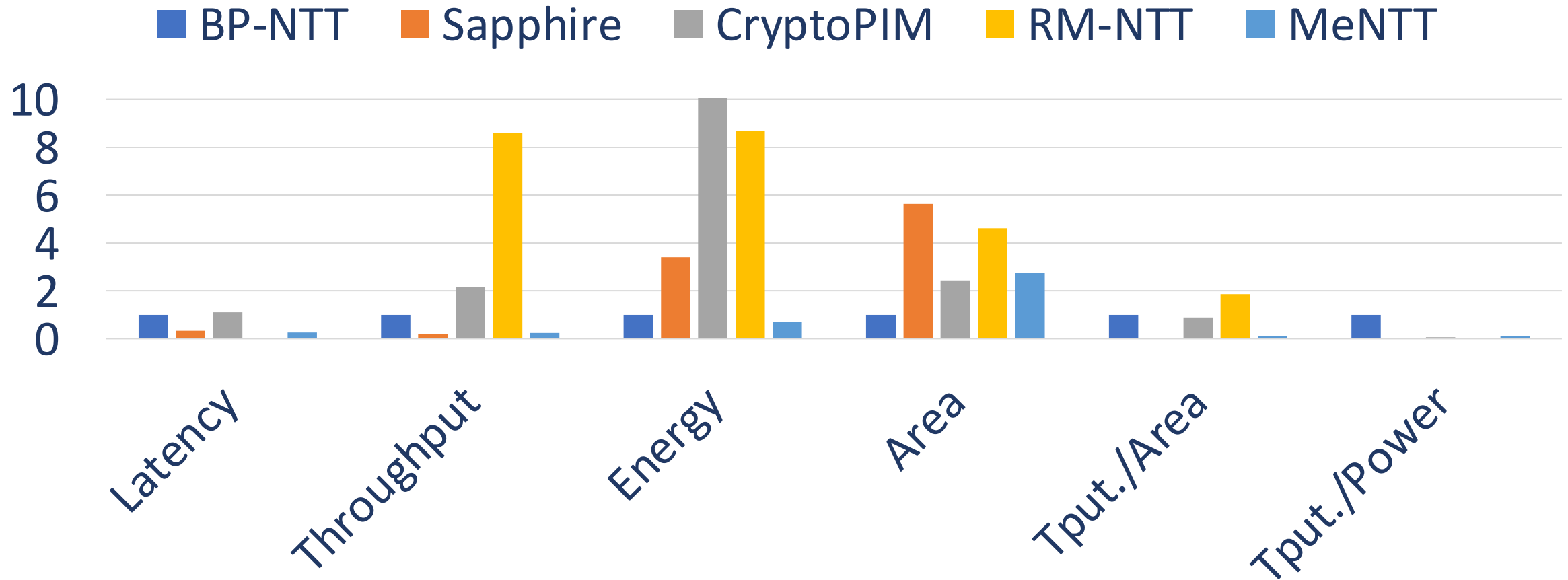
- Results are normalized to BP-NTT





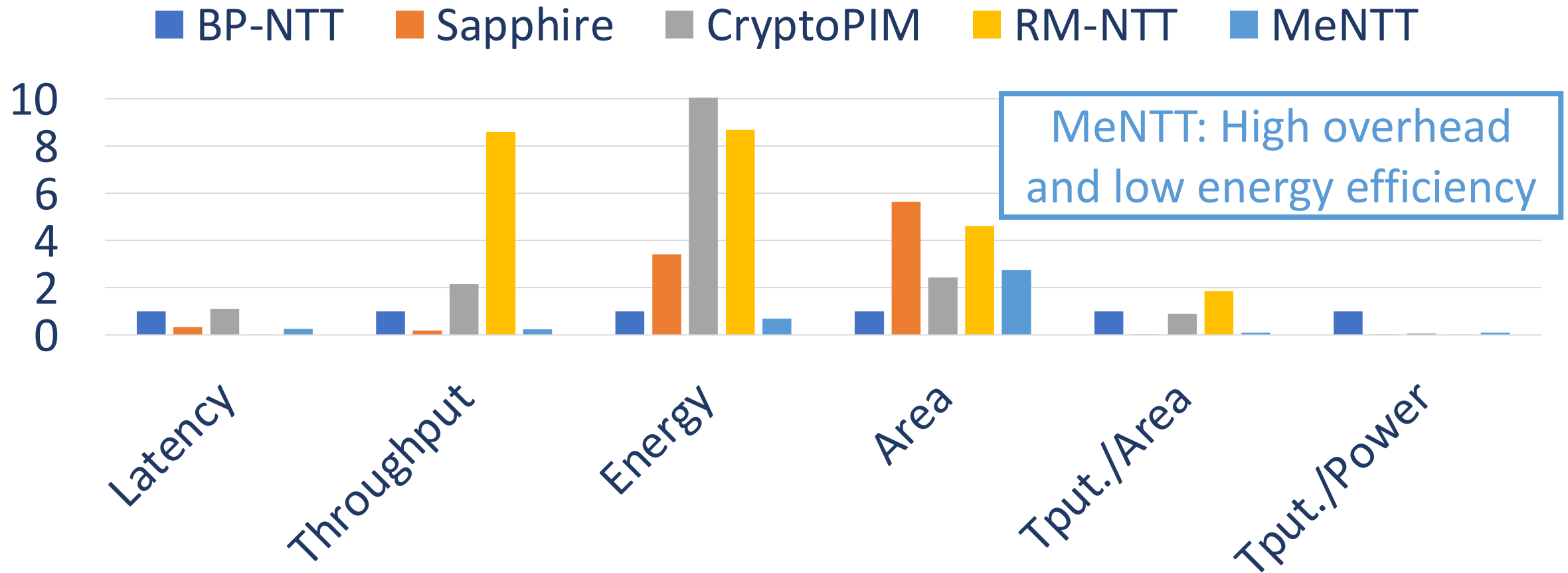
# Comparison among Designs

□ Results are normalized to BP-NTT



# Comparison among Designs

□ Results are normalized to BP-NTT



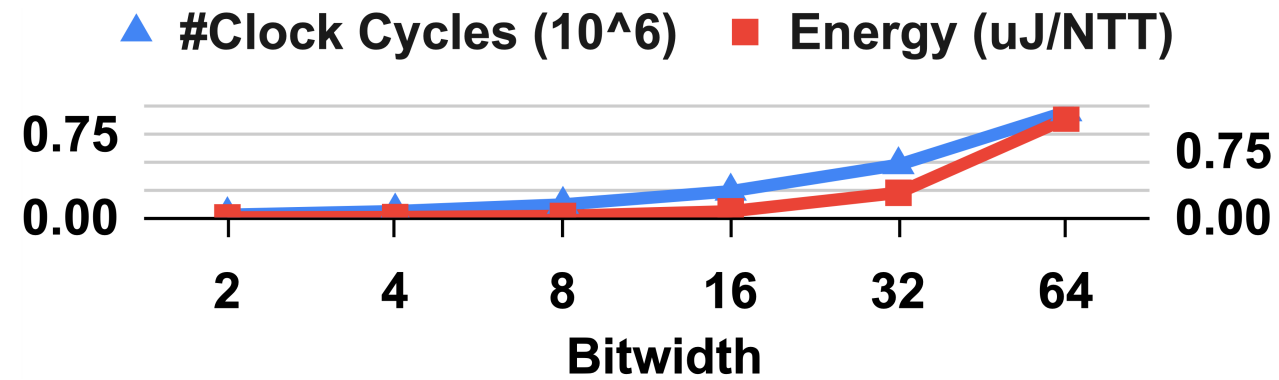
# Flexibility Analysis

---

# Flexibility Analysis

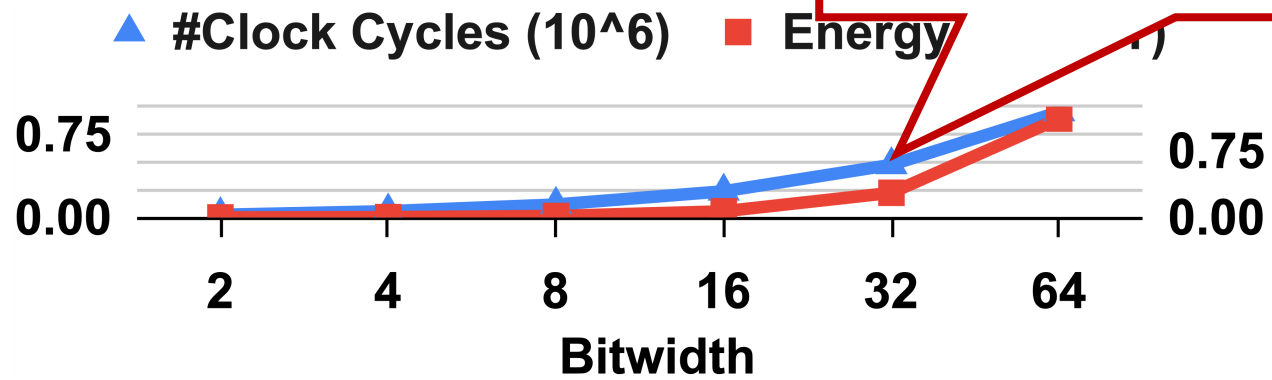
---

- With fixed polynomial order of 256



# Flexibility Analysis

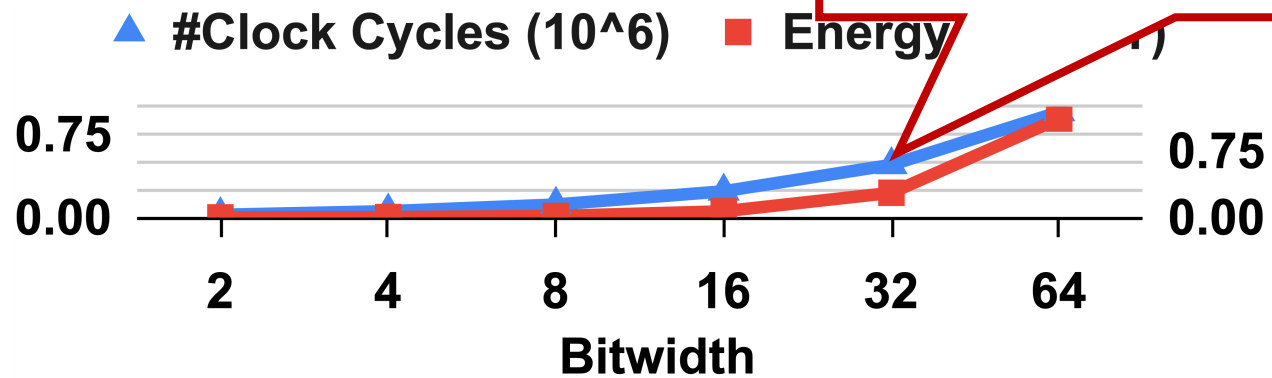
- With fixed polynomial order of 256



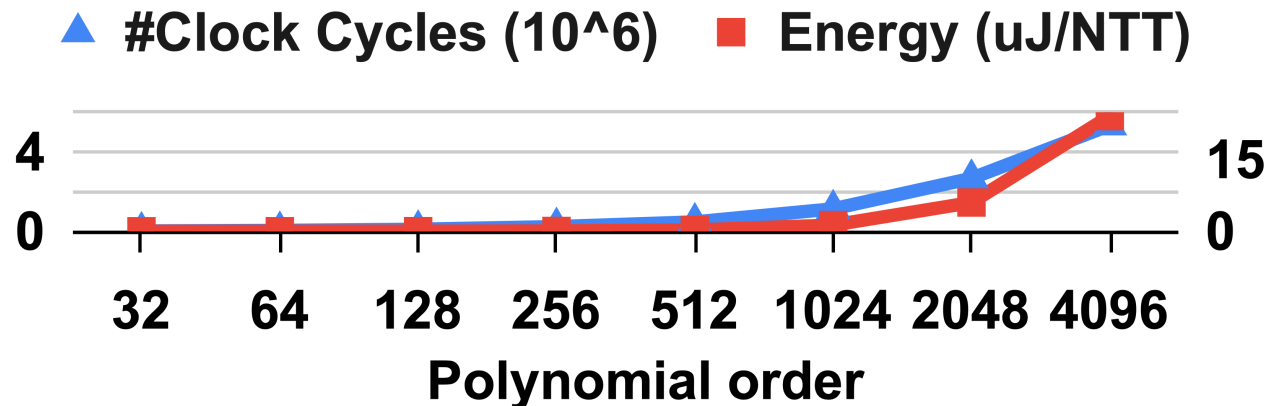
Parallelism decreases with increased bitwidth

# Flexibility Analysis

- With fixed polynomial order of 256

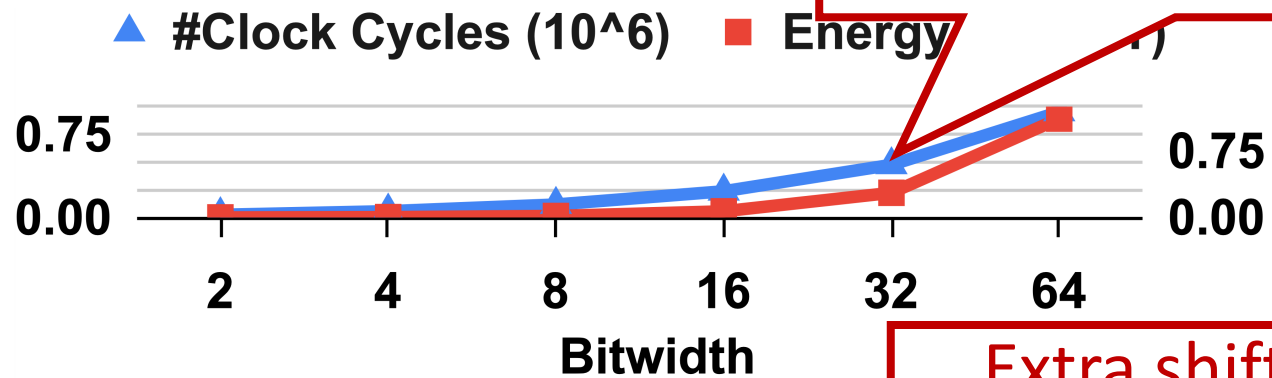


- With fixed bitwidth of 16

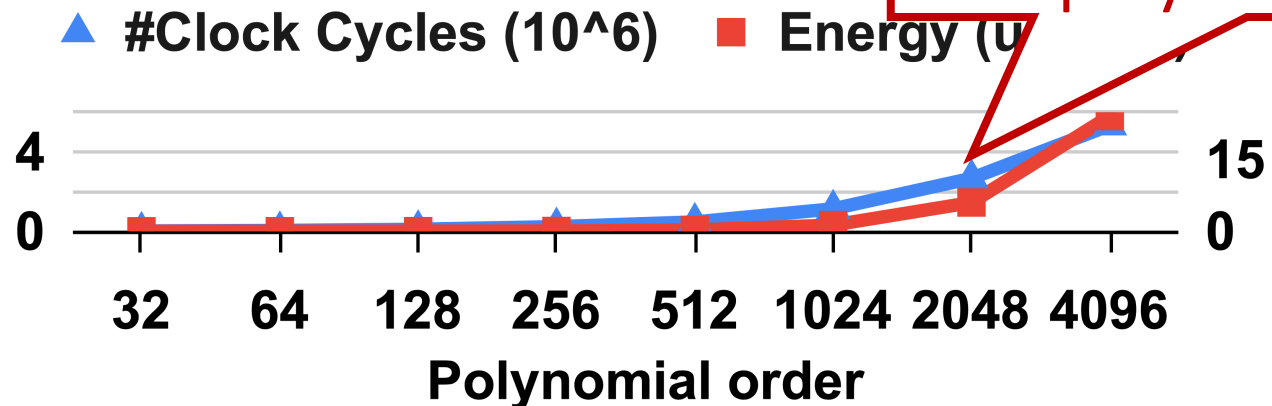


# Flexibility Analysis

- With fixed polynomial order of 256



- With fixed bitwidth of 16



# Conclusion

---



# Conclusion

---

***BP-NTT*** repurposes **LLC** to perform **bitline computing**



# Conclusion

---

***BP-NTT*** repurposes **LLC** to perform **bitline computing**



***BP-NTT*** uses **shift-optimized data alignment**



# Conclusion

---

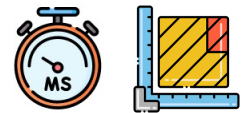
***BP-NTT*** repurposes **LLC** to perform **bitline computing**



***BP-NTT*** uses **shift-optimized data alignment**



***BP-NTT*** employs **bit-parallel modular multiplication**



# Conclusion

---

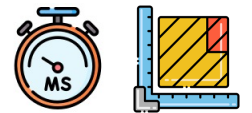
***BP-NTT*** repurposes **LLC** to perform **bitline computing**



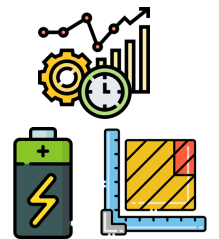
***BP-NTT*** uses **shift-optimized data alignment**



***BP-NTT*** employs **bit-parallel modular multiplication**



***BP-NTT*** can achieve **up to 138x** throughput-per-power than state-of-the-art with **minimal area**



# Q&A

---