

# Fast Erasure Coding for Data Storage: A Comprehensive Study of the Acceleration Techniques

*Tianli Zhou & Chao Tian  
Texas A&M University*



# Contents

- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- Proposed Coding Procedure and Evaluation
- Conclusion



# Contents

- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- Proposed Coding Procedure and Evaluation
- Conclusion



# Data Reliability



The image shows a grid of binary digits (0s and 1s) in green on a black background. Overlaid on this grid is the word "data" in a large, bold, blue sans-serif font.

```
01010100 01101000 01100101 00100000 01110001 01110101  
01101001 01100011 01101011 00100000 01100010 01110010  
01101111 01110111 01101110 00100000 01100110 01101111  
01111000 00100000 01101010 01110101 01101101 01100000  
01110011 00100000 01101111 01110110 01100101 01110010  
00100000 00110001 00110011 00100000 01101100 01100001  
01111010 01111001 00100000 01101000 01101111 01100111  
01010100 01101000 00101010 01100000 01100011 01101011  
01101001 01100001 10101011 01100000 01100010 01110010  
01101111 01110111 01101110 00100000 01100110 01101111  
01111000 00100000 01101010 01110101 01101101 01110000  
01110011 00100000 01101111 01110110 01100101 01110010  
00100000 00110001 00110011 00100000 01101100 01100001  
01111010 01111001 00100000 01101000 01101111 01100111
```



# Data Reliability

- Disaster happens

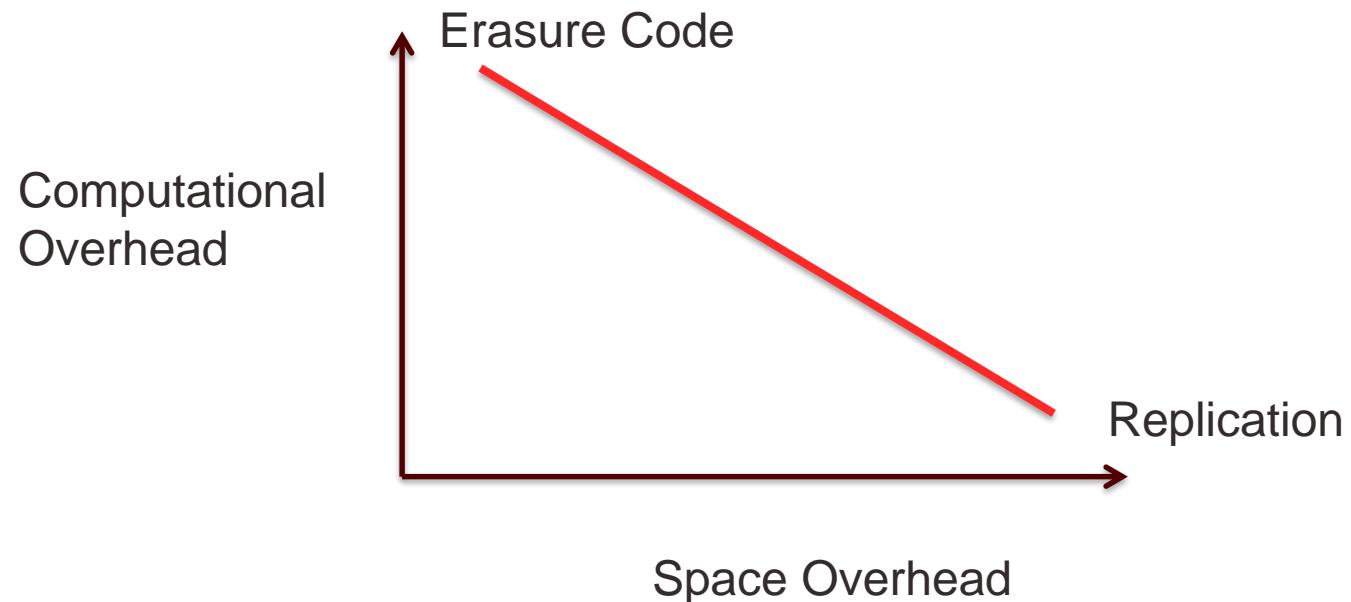


# Data Reliability

- Disaster happens
  - Disk failure
  - Server malfunction
  - and more....

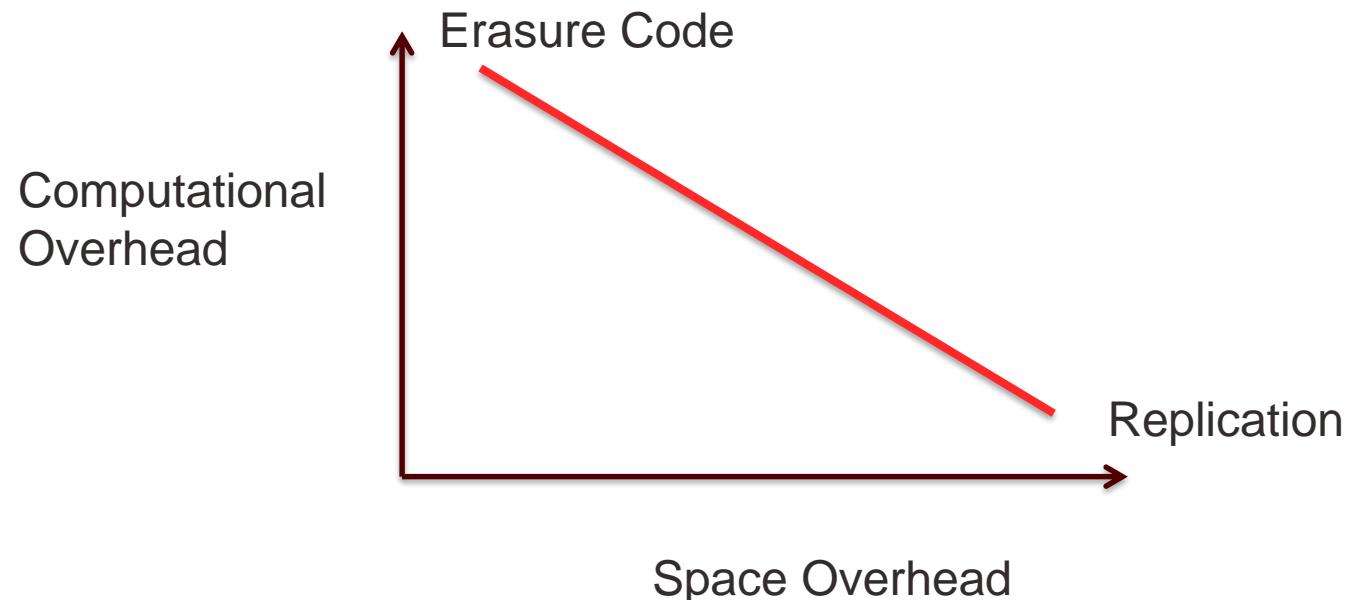


# Replication - Erasure Code



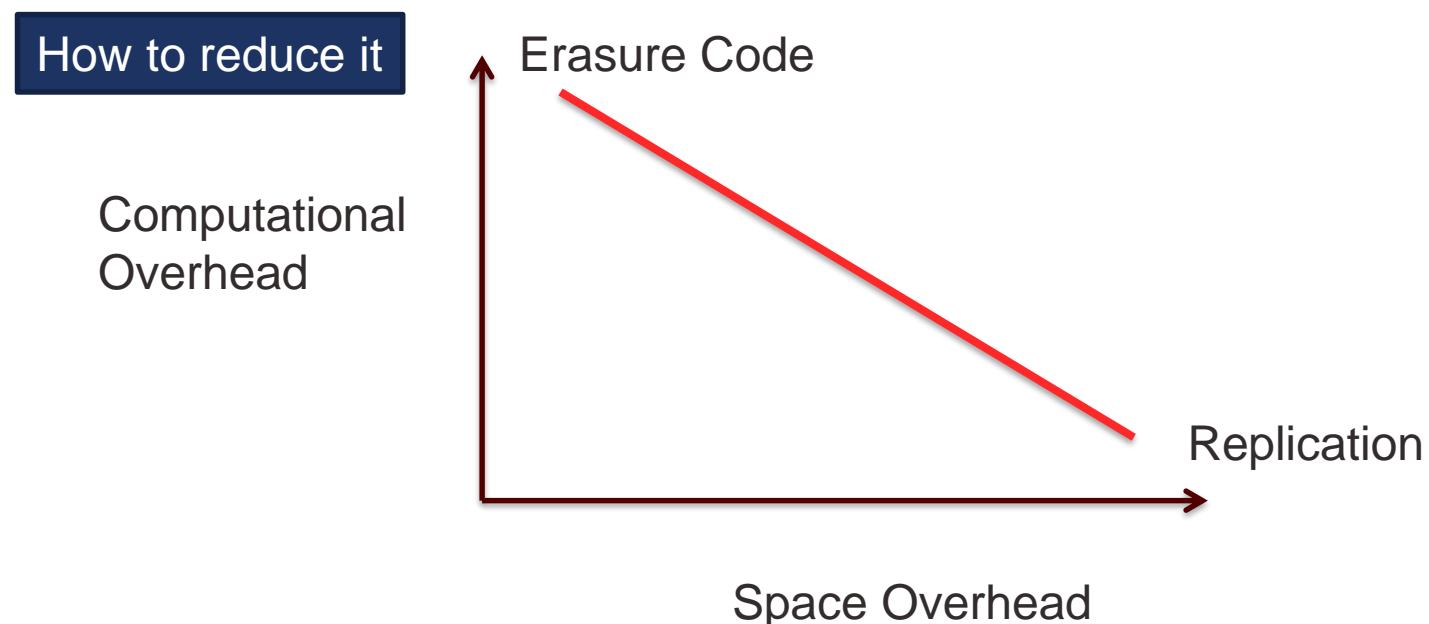
# Replication - Erasure Code

- Minimum space overhead
- Much more computational load
- Erasure code relies on Galois Field (finite field) ops



# Replication - Erasure Code

- Minimum space overhead
- Much more computational load
- Erasure code relies on Galois Field (finite field) ops



# Acceleration Techniques

Optimize  
Coding Matrix

Schedule  
Computing  
Order

Covert to XOR  
operations

Apply SIMD

Specially  
Designed Code

Optimize for  
CPU Cache



# Acceleration Techniques

---

Fast GF  
Coding

Covert to XOR  
operations

Optimize  
Coding Matrix

Schedule  
Computing  
Order

---

Coding  
Theory

---

Specially  
Designed Code

---

Computation  
Resource

---

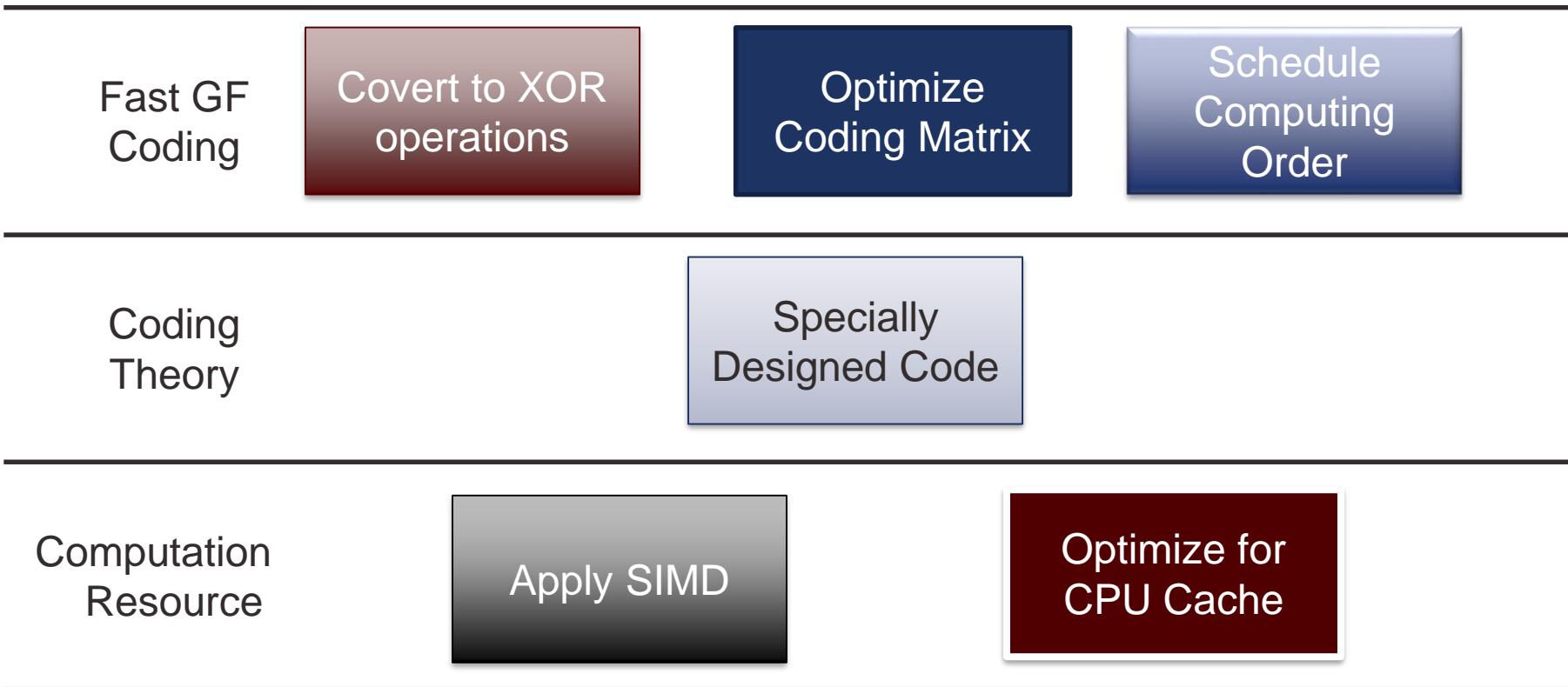
Apply SIMD

Optimize for  
CPU Cache

---



# Acceleration Techniques



?

# Question to Answer

- Most effective technique(s)?
- Utilize together?
- Components to be optimized?
- Problem is still important
  - Reduce energy consumption
  - Virtualized environment
  - Mobile / embedded system
  - Future I/O technology

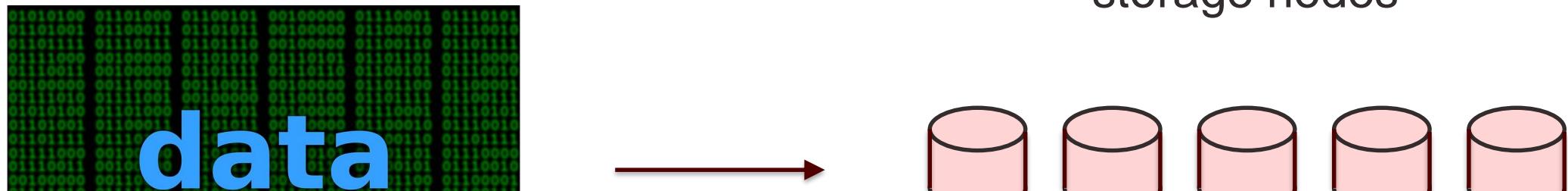


# Contents

- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- Proposed Coding Procedure and Evaluation
- Conclusion

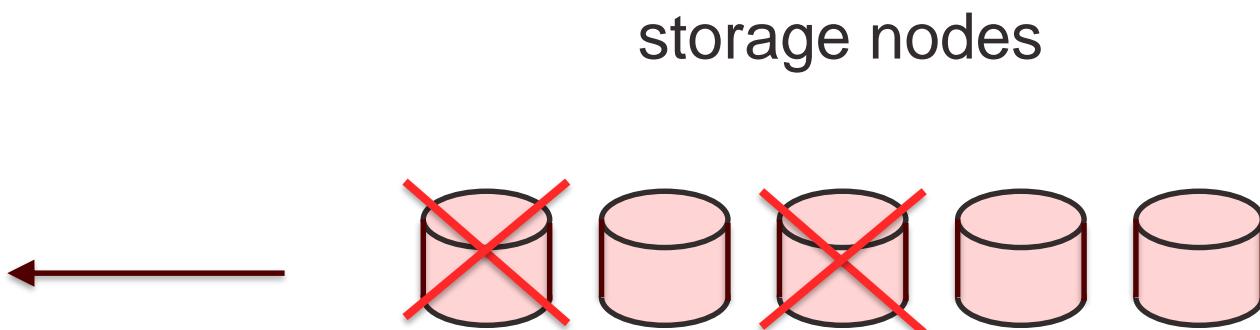


# Erasures Code



Store coded data to multiple nodes

# Erasasure Code



Recover even nodes failures happen



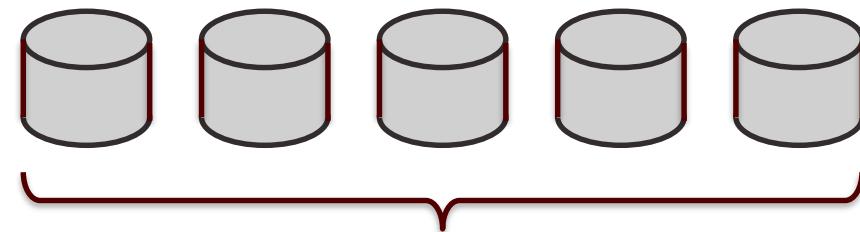
# Erasures Code



$k=3$   
segments

Encode  
→

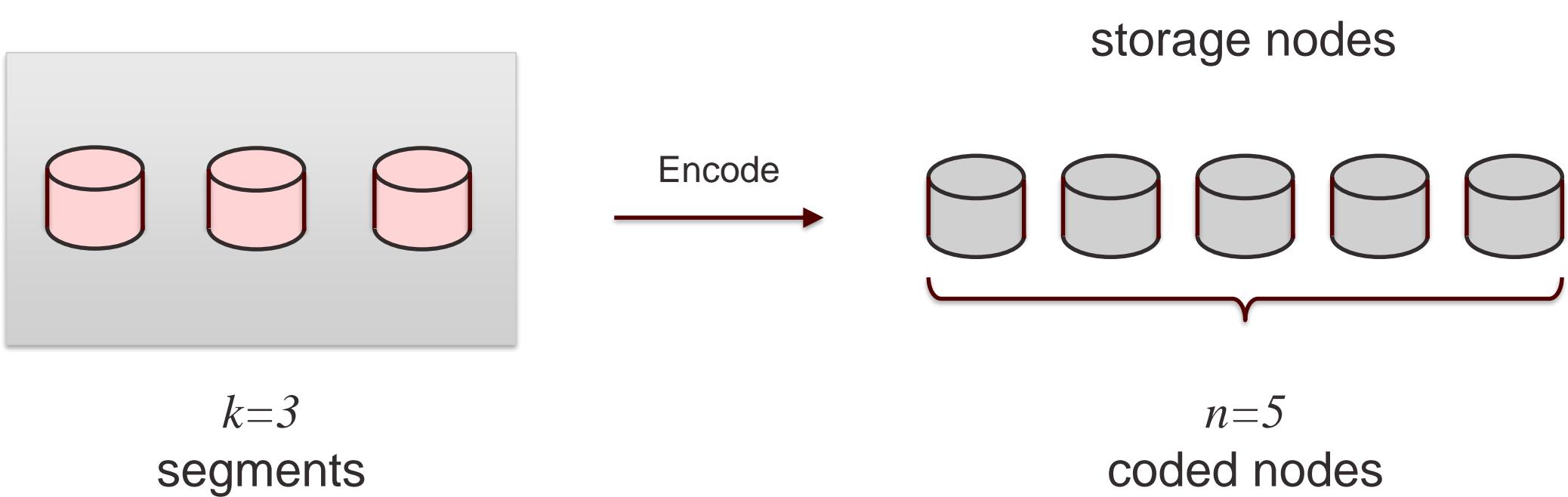
storage nodes



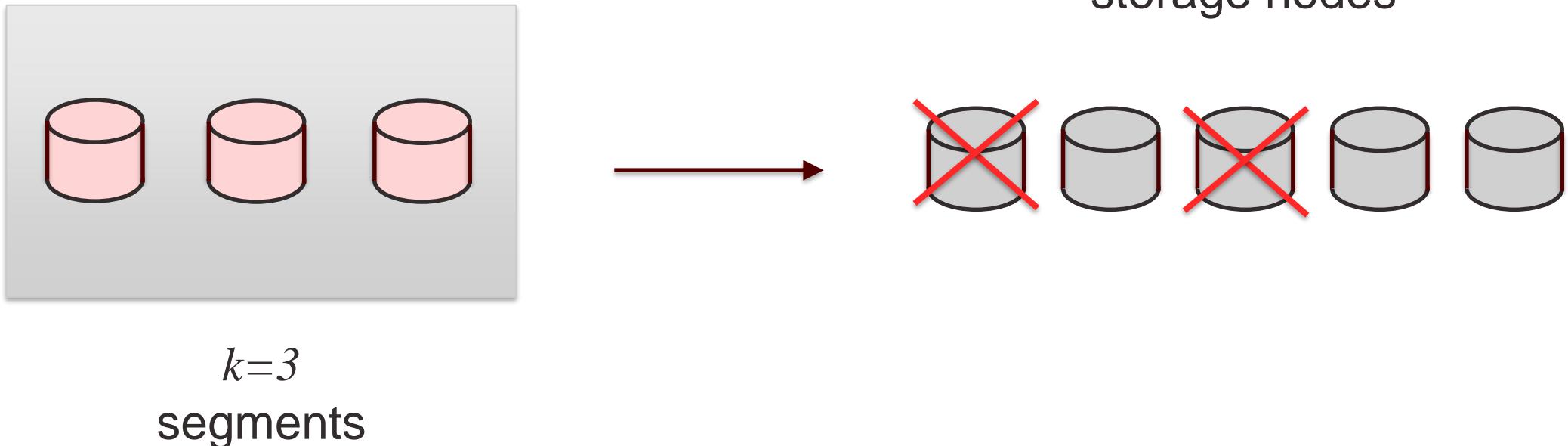
$n=5$   
coded nodes



# Erasasure Code



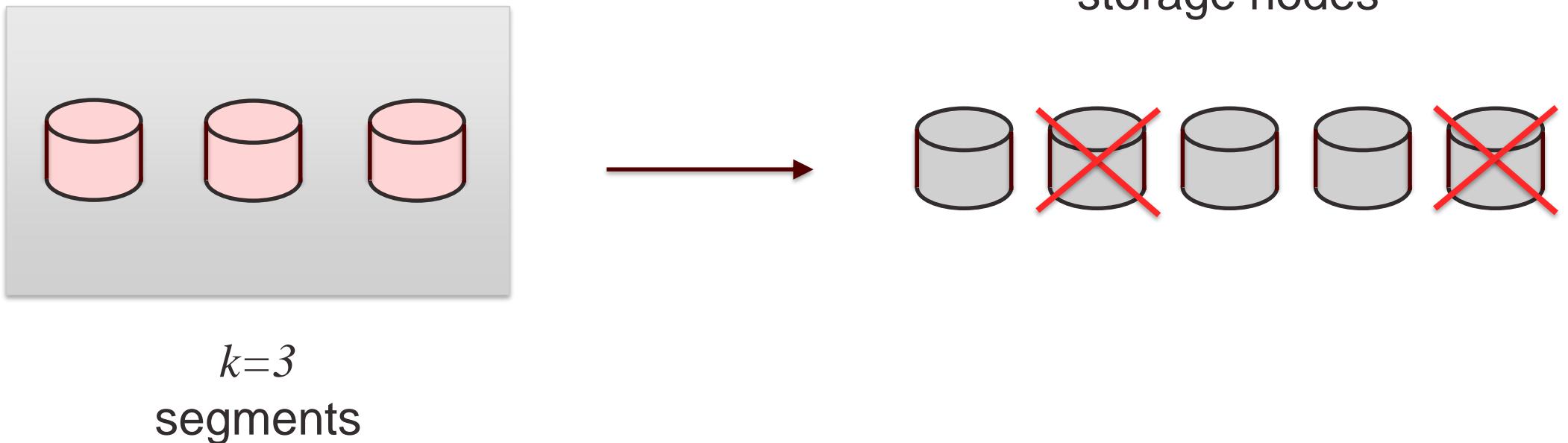
# Erasasure Code



MDS code: recover data  
from any  $k$  surviving nodes



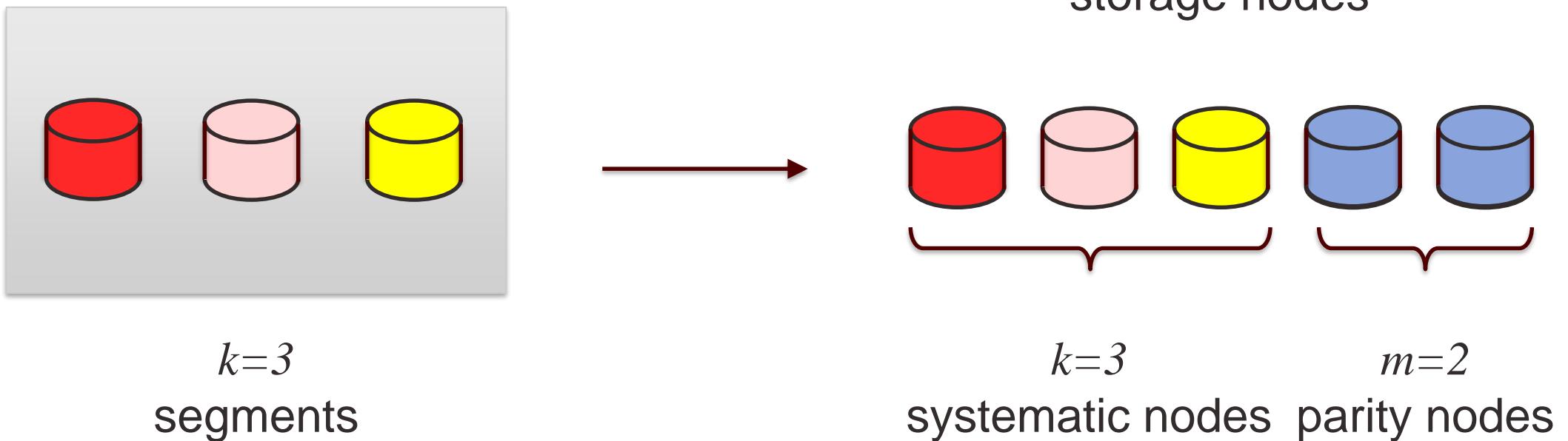
# Erasasure Code



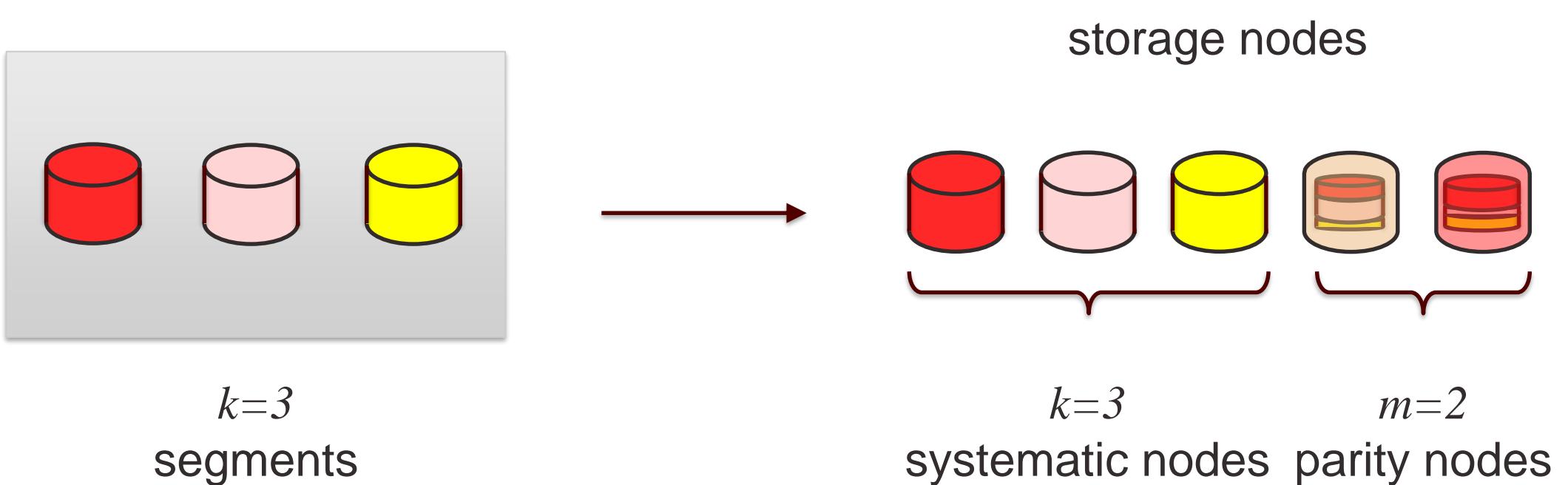
MDS code: recover data  
from any  $k$  surviving nodes



# Erasasure Code



# Erasasure Code



# Encoding

- Matrix multiplication on finite field

$$\begin{array}{c} \text{red cylinder} \\ \text{pink cylinder} \\ \text{yellow cylinder} \end{array} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} \\ p_{2,0} & p_{2,1} \end{bmatrix} = \begin{array}{c} \text{red cylinder} \\ \text{pink cylinder} \\ \text{yellow cylinder} \\ \text{brown cylinder} \\ \text{red cylinder} \end{array}$$

# Encoding

- Matrix multiplication on finite field

$$\begin{array}{c} \text{red cylinder} \\ \text{pink cylinder} \\ \text{yellow cylinder} \end{array} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} \\ p_{2,0} & p_{2,1} \end{bmatrix} = \begin{array}{c} \text{red cylinder} \\ \text{pink cylinder} \\ \text{yellow cylinder} \\ \text{brown cylinder} \\ \text{red cylinder} \end{array}$$

All elements and ops in  $GF(2^w)$

Galois Field (GF)  
Finite field



# Decoding

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} & p_{1,1} \\ p_{2,0} & & p_{2,1} \end{bmatrix}$$



# Decoding

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} & p_{1,1} \\ p_{2,0} & & p_{2,1} \end{bmatrix}$$



# Decoding

- Submatrix corresponding to surviving data

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} & p_{1,1} \\ p_{2,0} & & p_{2,1} \end{bmatrix}$$



$$\begin{array}{c|cc} 0 & p_{0,0} & p_{0,1} \\ 0 & p_{1,0} & p_{1,1} \\ 1 & p_{2,0} & p_{2,1} \end{array}$$

# Decoding

- Submatrix corresponding to surviving data
- Invert

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} & p_{1,1} \\ p_{2,0} & & p_{2,1} \end{bmatrix}$$



$$\begin{array}{c|cc|c} 0 & p_{0,0} & p_{0,1} & -1 \\ 0 & p_{1,0} & p_{1,1} & \\ 1 & p_{2,0} & p_{2,1} & \end{array}$$

# Decoding

- Submatrix corresponding to surviving data
- Invert
- Multiplication

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \mathbf{I} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} p_{0,0} & p_{0,1} \\ p_{1,0} & \mathbf{P} & p_{1,1} \\ p_{2,0} & & p_{2,1} \end{bmatrix}$$

$$\begin{array}{c} \text{Yellow cylinder} \\ \text{Brown cylinder} \\ \text{Red cylinder} \end{array} \times \begin{array}{c} \begin{array}{c|cc} 0 & p_{0,0} & p_{0,1} \\ 0 & p_{1,0} & p_{1,1} \\ 1 & p_{2,0} & p_{2,1} \end{array}^{-1} \\ = \end{array} \begin{array}{c} \text{Red cylinder} \\ \text{Pink cylinder} \\ \text{Yellow cylinder} \end{array}$$



# Cauchy Reed-Solomon Codes

- Direct assign  $P$  to Cauchy matrix

$$\left[ \begin{array}{cccc} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{array} \middle| \begin{array}{ccc} p_{0,0} & \cdots & p_{0,m-1} \\ p_{1,0} & \cdots & p_{1,m-1} \\ \vdots & \ddots & \vdots \\ p_{k-1,0} & \cdots & p_{k-1,m-1} \end{array} \right]$$



# Cauchy Reed-Solomon Codes

$$X = (x_1, \dots, x_k) \quad Y = (y_1, \dots, y_m)$$

$$C = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \frac{x_1 + y_1}{1} & \frac{x_1 + y_2}{1} & \dots & \frac{x_1 + y_m}{1} \\ \frac{x_2 + y_1}{1} & \frac{x_2 + y_2}{1} & \dots & \frac{x_2 + y_m}{1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \\ \frac{x_k + y_1}{1} & \frac{x_k + y_2}{1} & \dots & \frac{x_k + y_m}{1} \end{bmatrix}$$

# Cauchy Reed-Solomon Codes

$$X = (x_1, \dots, x_k) \quad Y = (y_1, \dots, y_m)$$

$$C = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \frac{x_1 + y_1}{1} & \frac{x_1 + y_2}{1} & \dots & \frac{x_1 + y_m}{1} \\ \frac{x_2 + y_1}{1} & \frac{x_2 + y_2}{1} & \dots & \frac{x_2 + y_m}{1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \\ \frac{x_k + y_1}{1} & \frac{x_k + y_2}{1} & \dots & \frac{x_k + y_m}{1} \end{bmatrix}$$

All elements and ops in  $GF(2^w)$



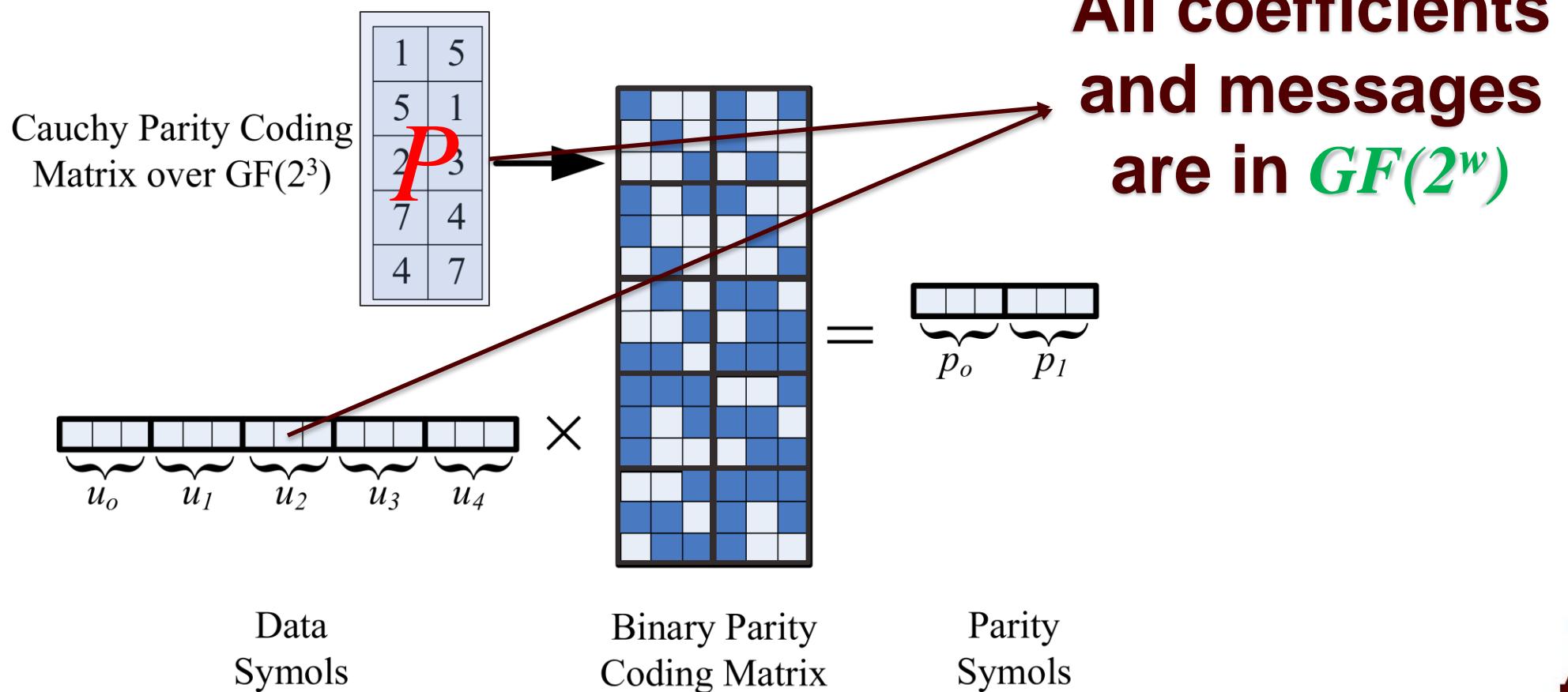
# Cauchy Reed-Solomon Codes

- Direct assign  $P$  to Cauchy matrix

$$\begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \left[ \begin{array}{c|ccc} \frac{1}{x_1 + y_1} & \frac{1}{x_1 + y_2} & \cdots & \frac{1}{x_1 + y_m} \\ \hline \frac{1}{x_2 + y_1} & \frac{1}{x_2 + y_2} & \cdots & \frac{1}{x_2 + y_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{x_k + y_1} & \frac{1}{x_k + y_2} & \cdots & \frac{1}{x_k + y_m} \end{array} \right] P$$

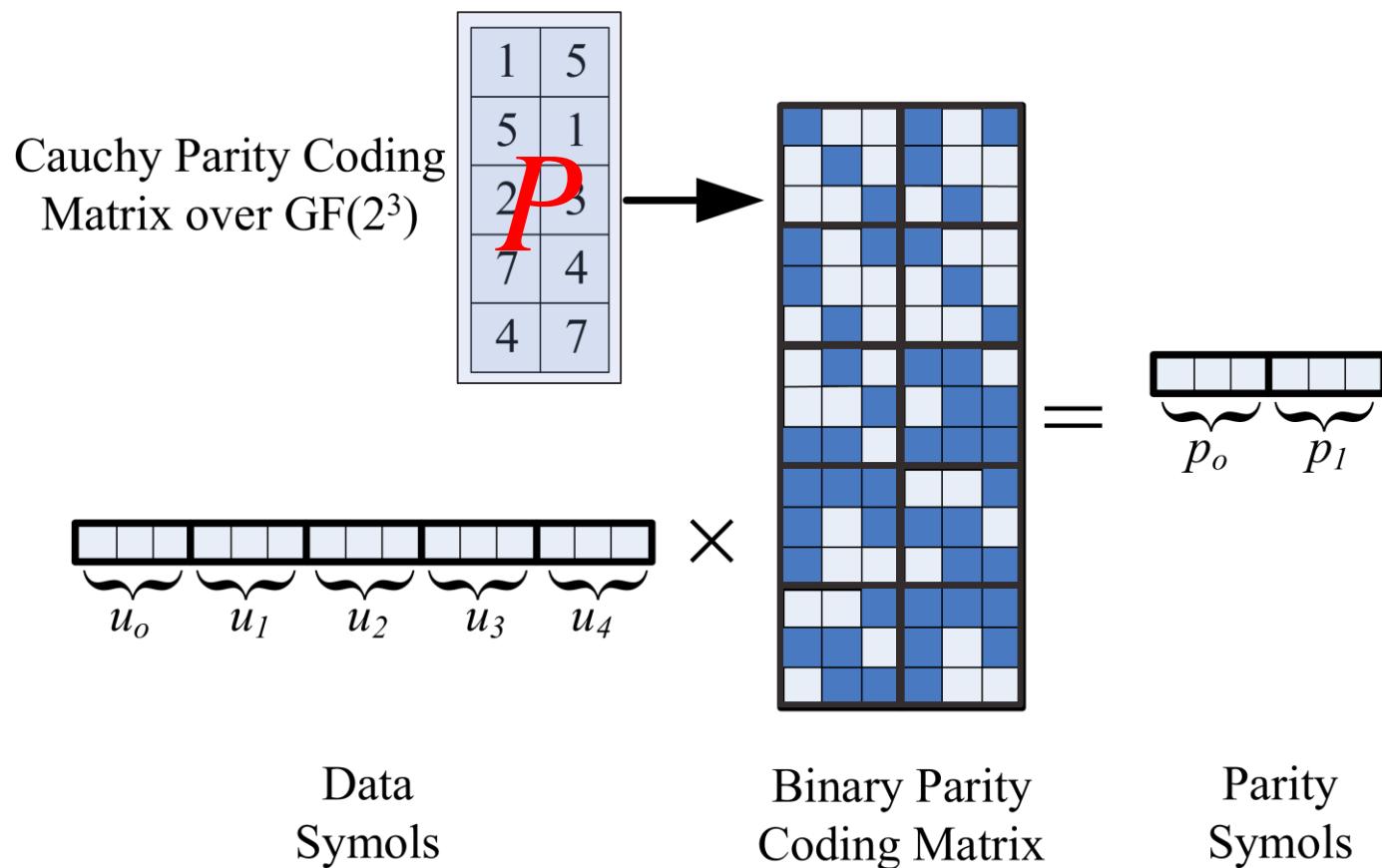
# Fast GF Ops – Binary Representation

- Binary representation of GF elements



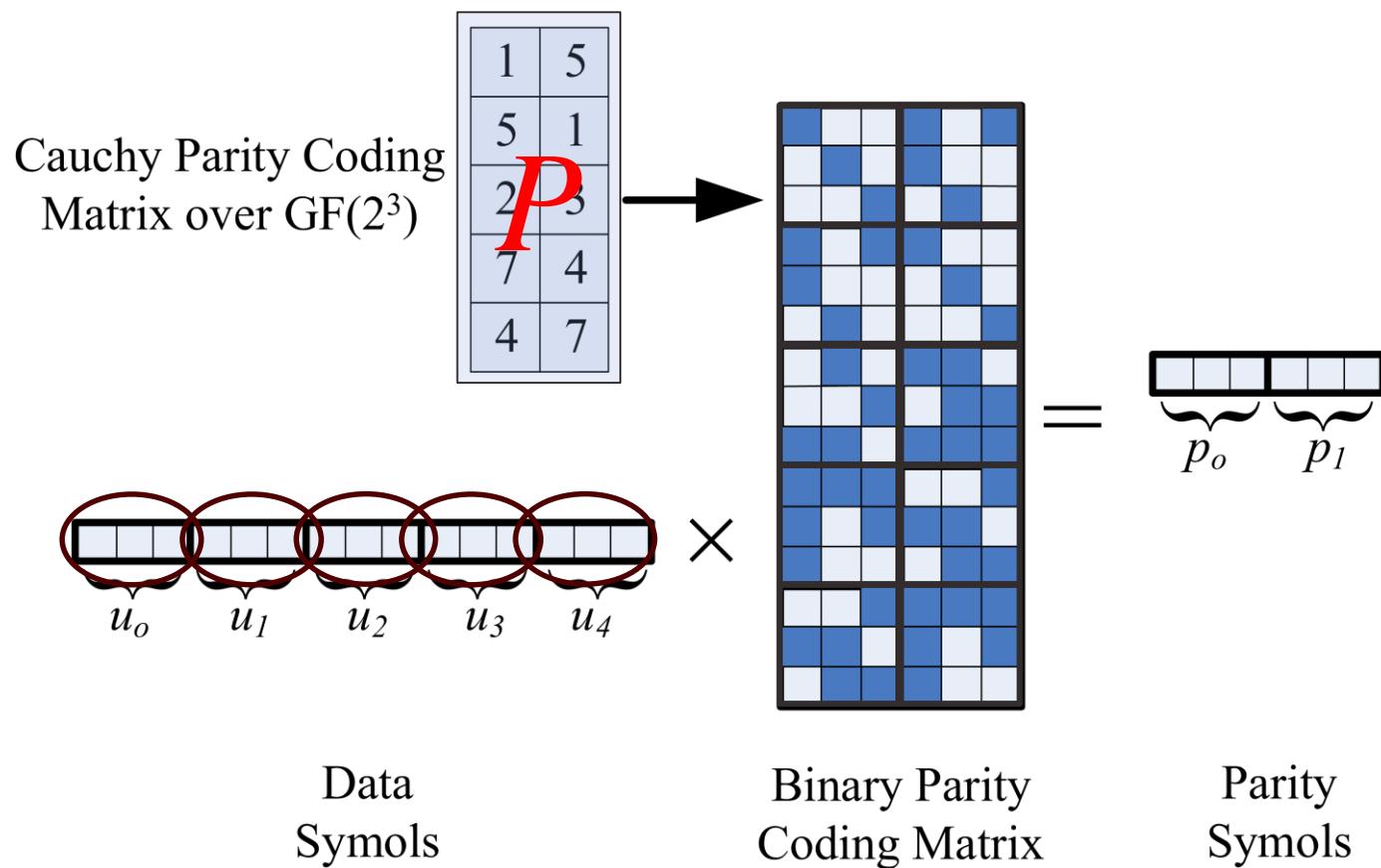
# Fast GF Ops – Binary Representation

- Message bit-vector representation



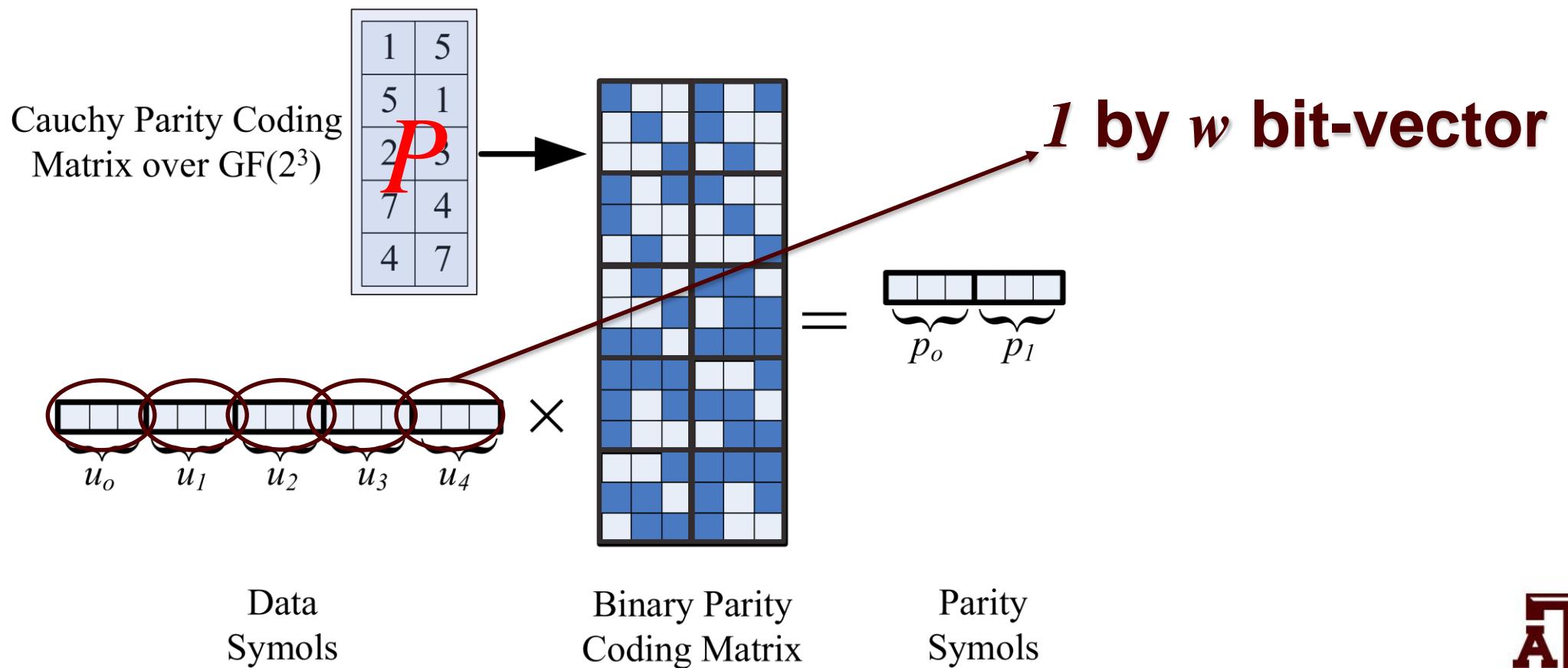
# Fast GF Ops – Binary Representation

- Message bit-vector representation



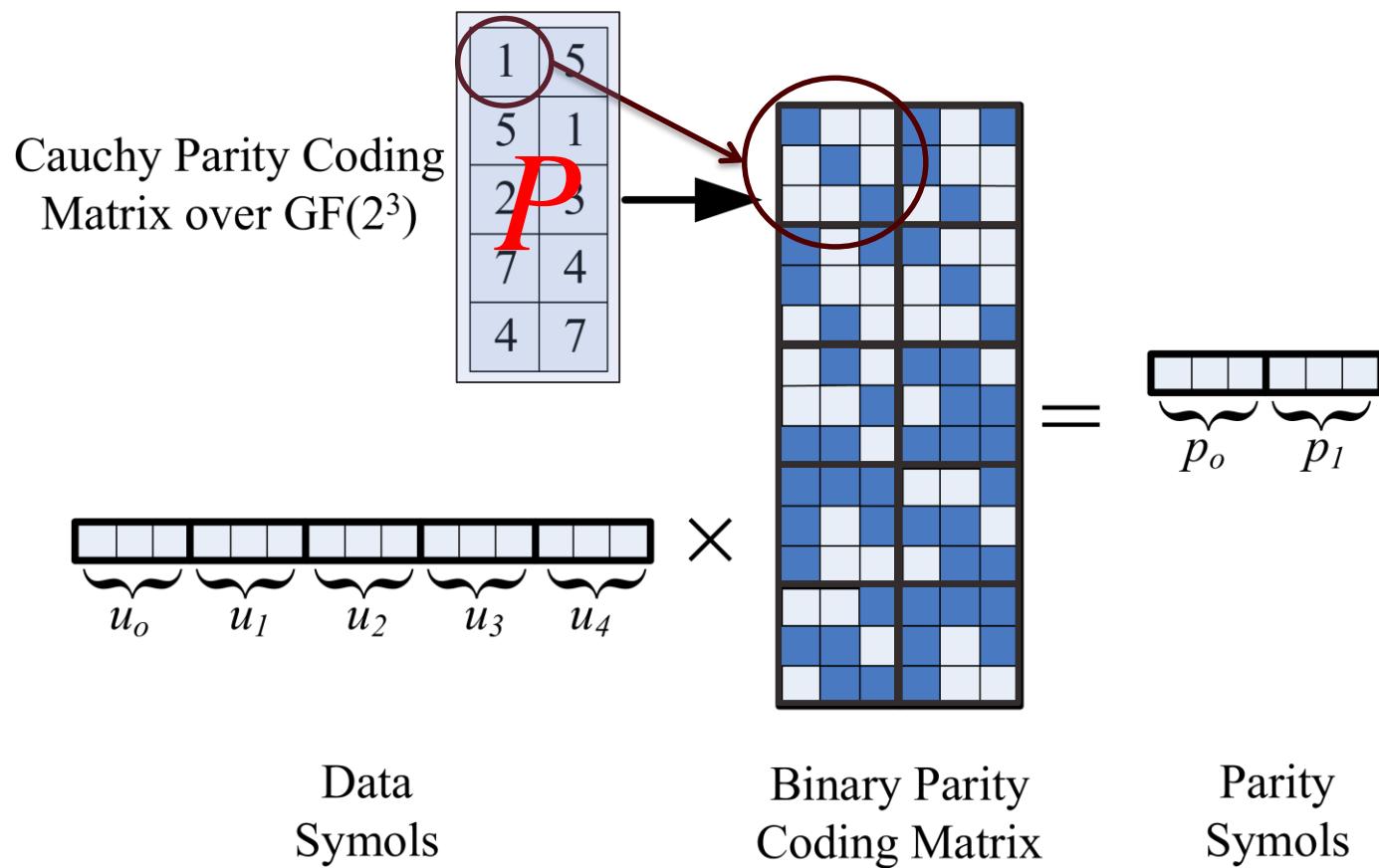
# Fast GF Ops – Binary Representation

- Message bit-vector representation



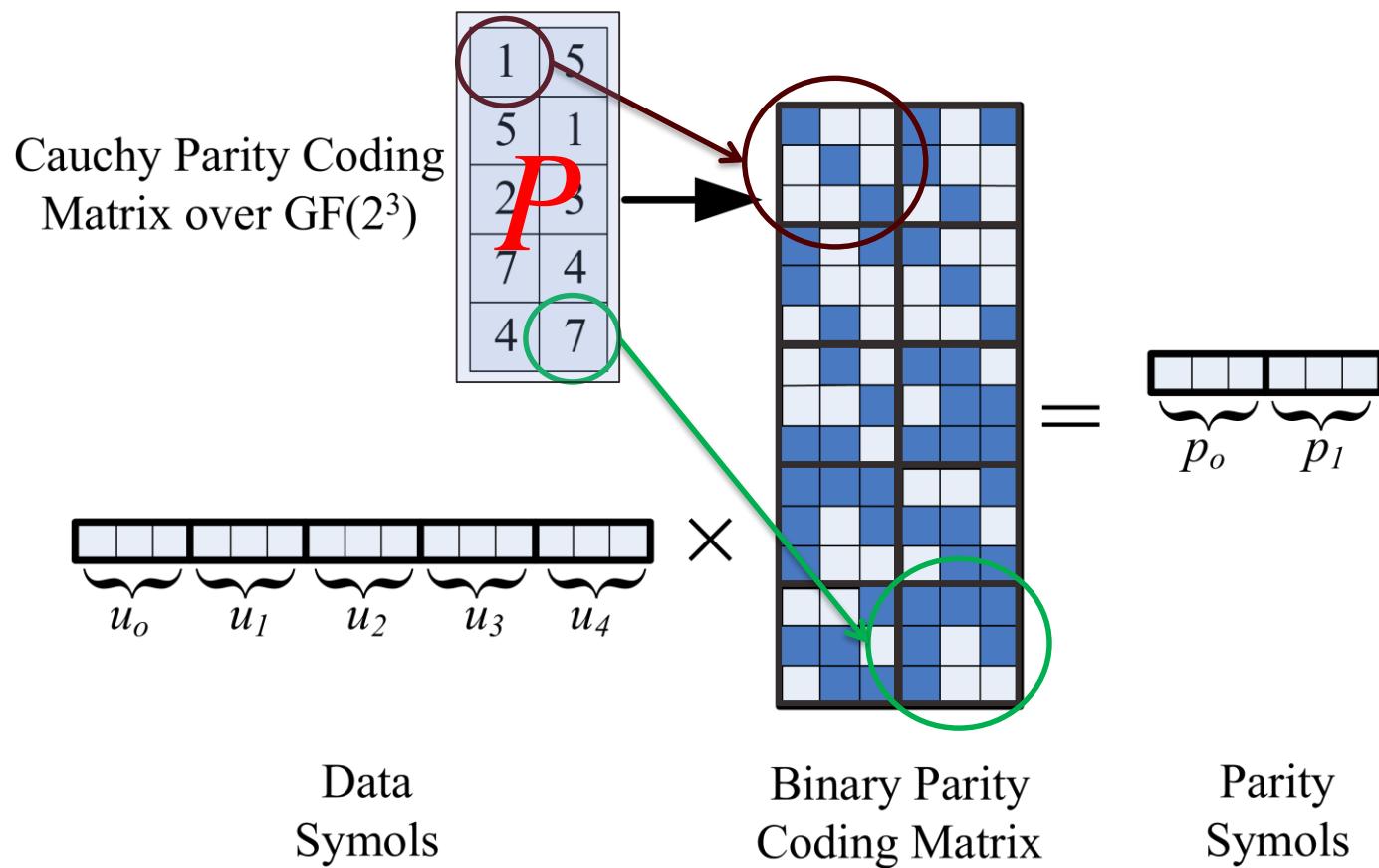
# Fast GF Ops – Binary Representation

- Bitmatrix representation



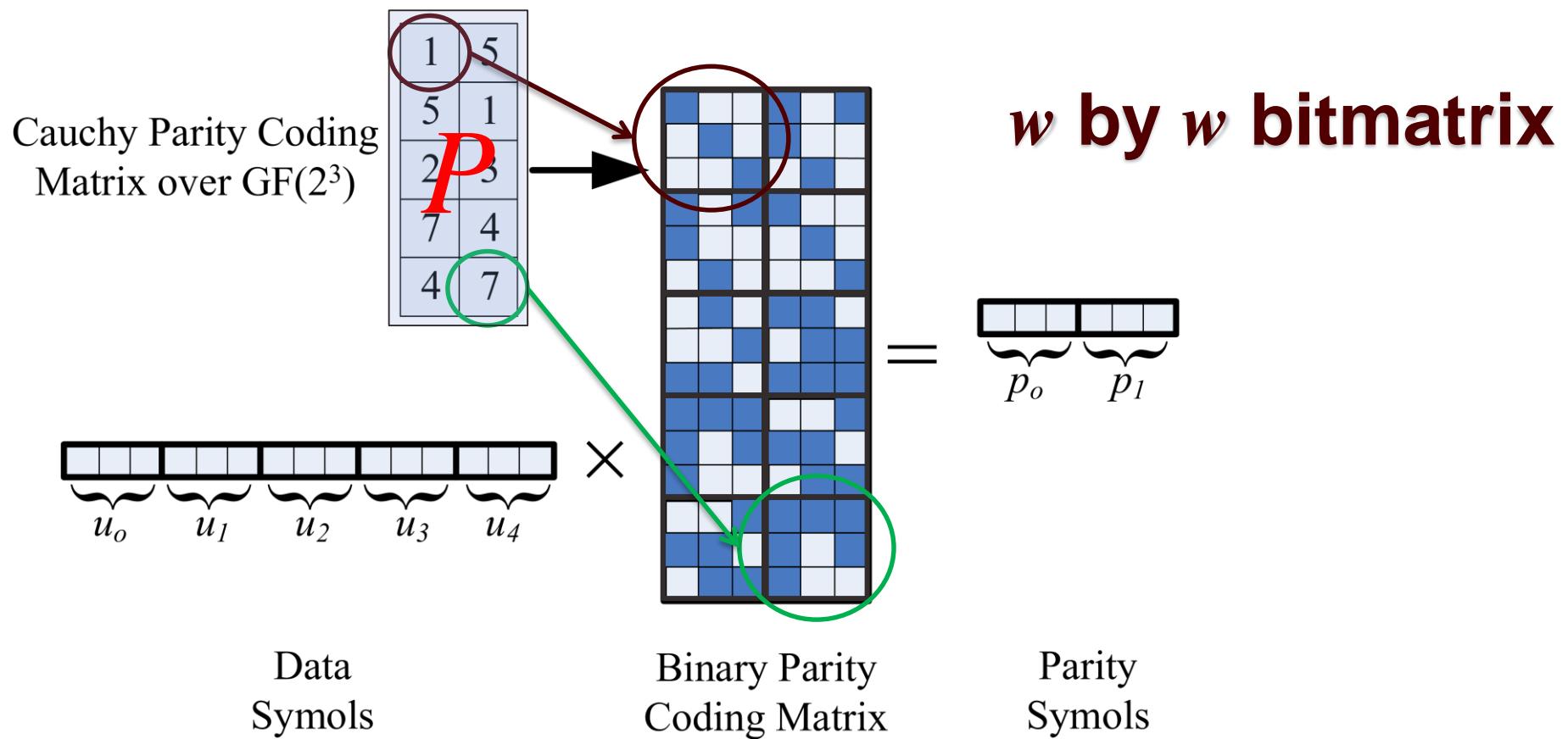
# Fast GF Ops – Binary Representation

- Bitmatrix representation



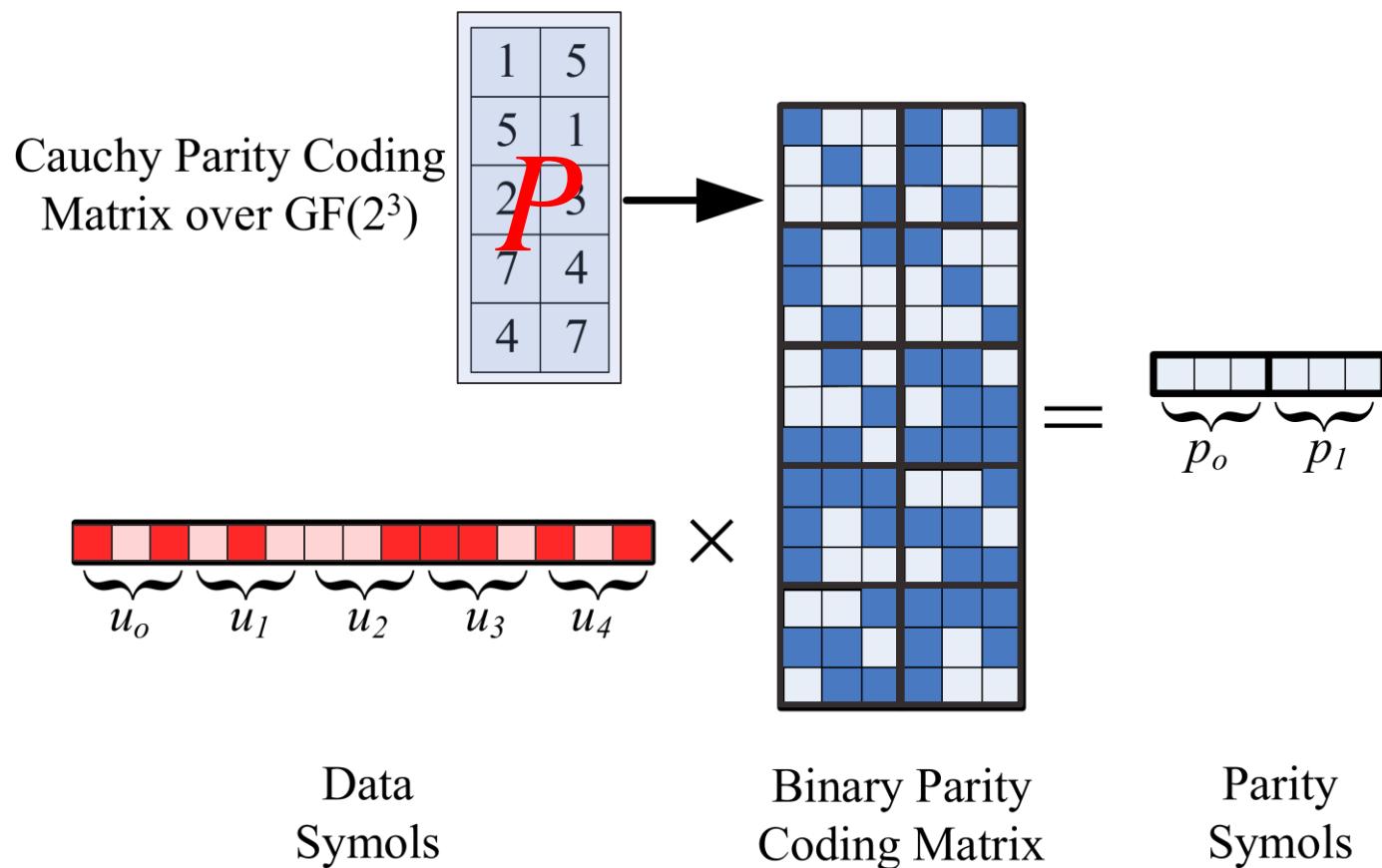
# Fast GF Ops – Binary Representation

- Bitmatrix representation



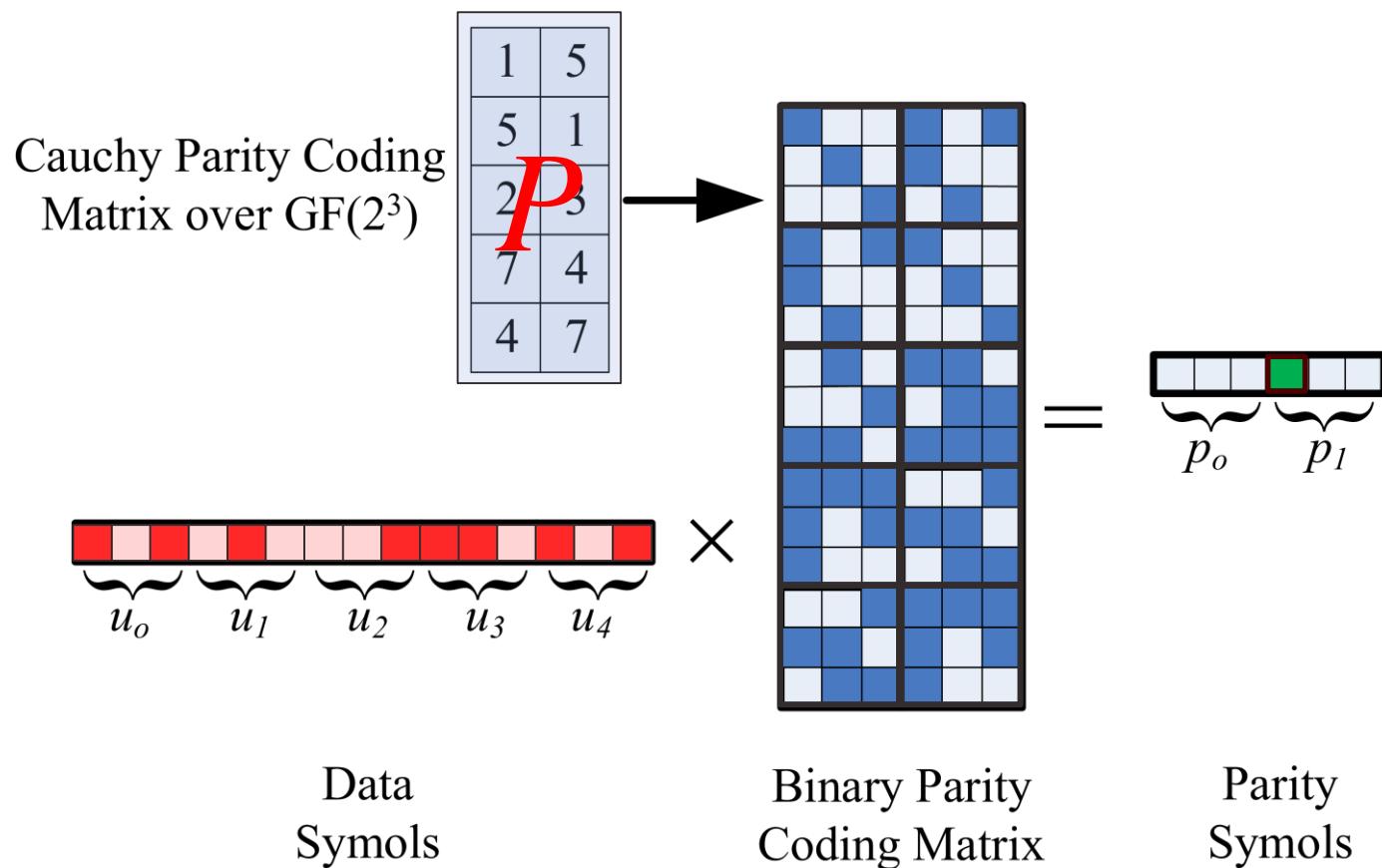
# Fast GF Ops – Binary Representation

- Matrix multiplication on GF -> XORs



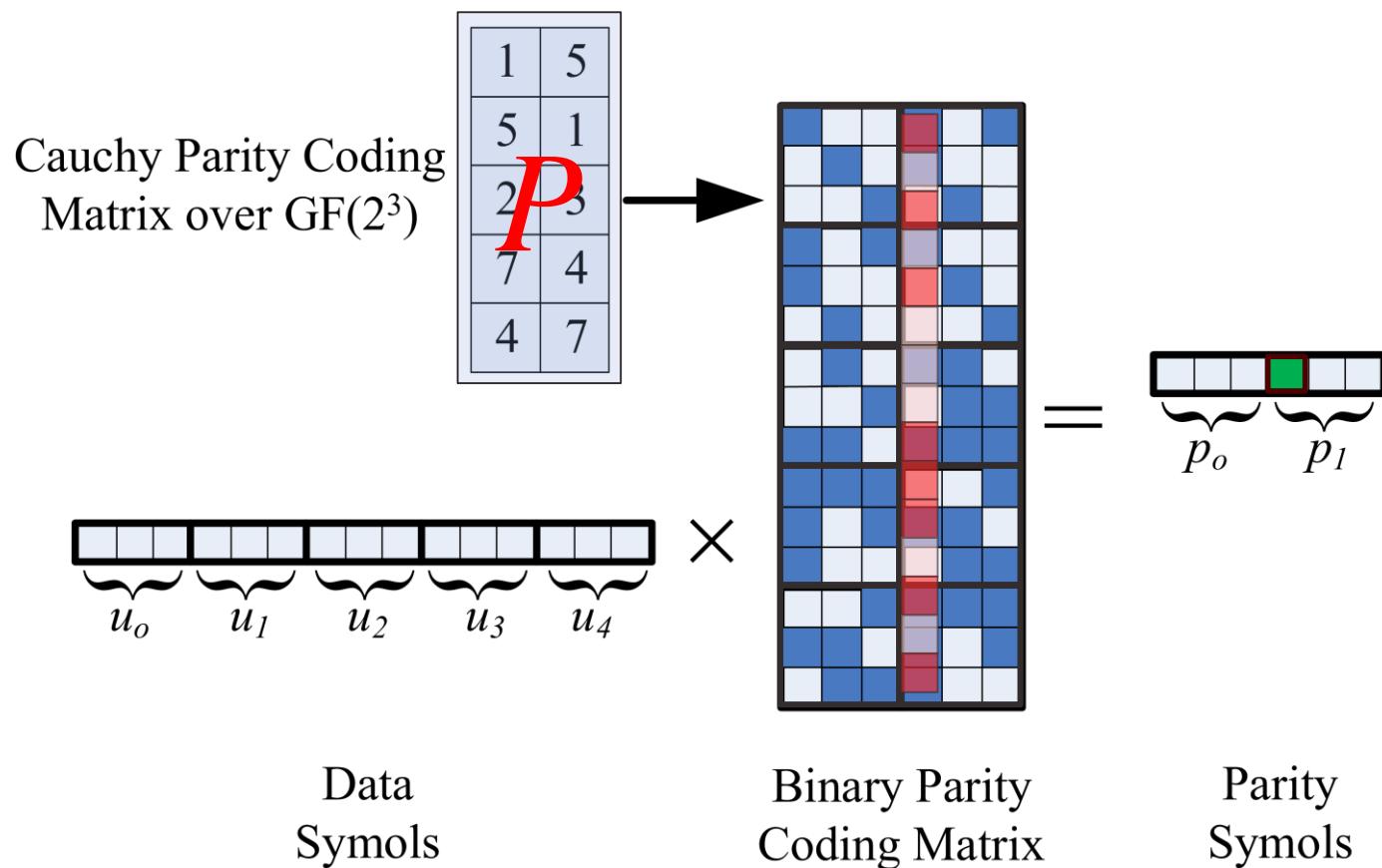
# Fast GF Ops – Binary Representation

- Matrix multiplication on GF -> XORs



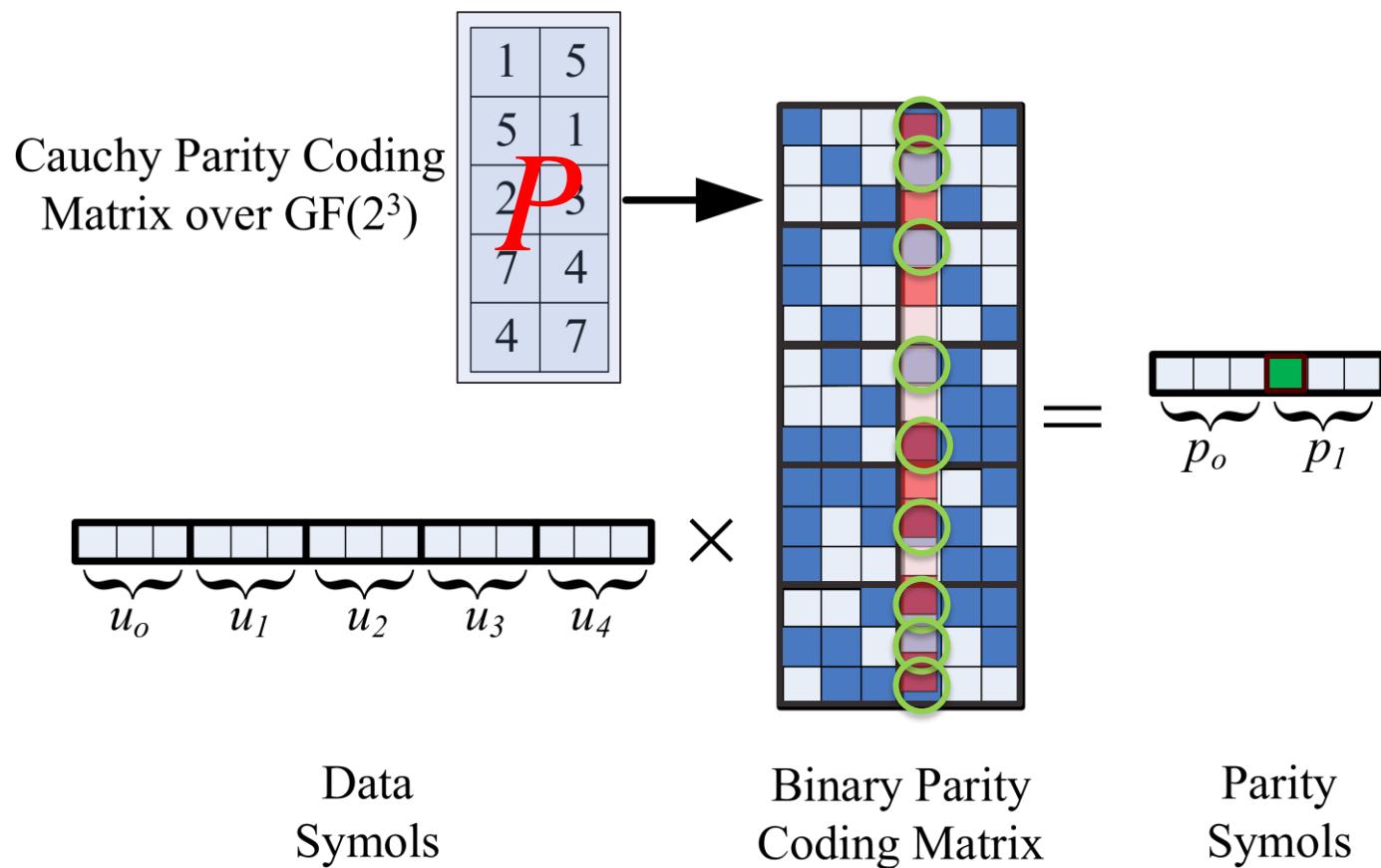
# Fast GF Ops – Binary Representation

- Matrix multiplication on GF -> XORs



# Fast GF Ops – Binary Representation

- Matrix multiplication on GF -> XORs

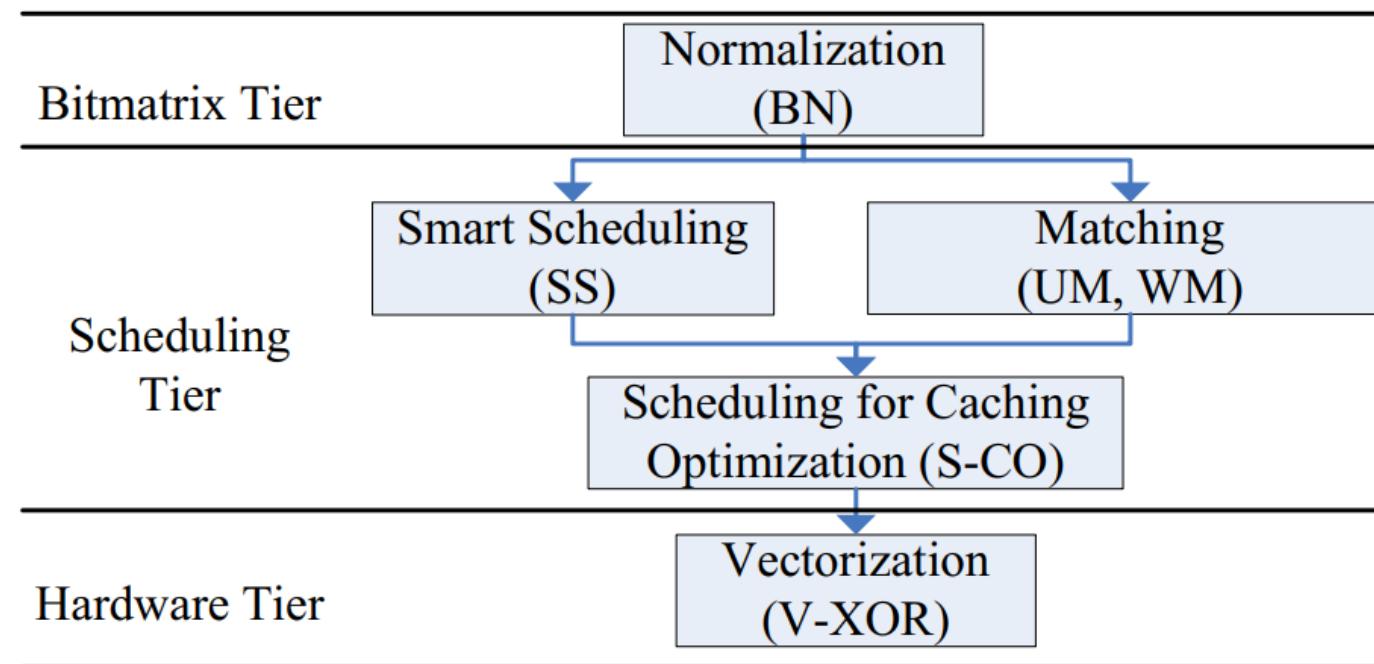


# Contents

- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- Proposed Coding Procedure and Evaluation
- Conclusion

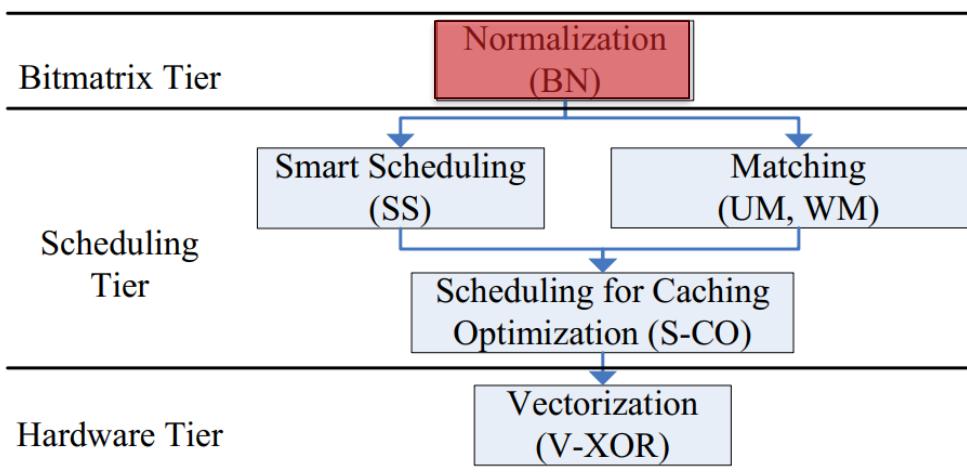


# Techniques Tiers



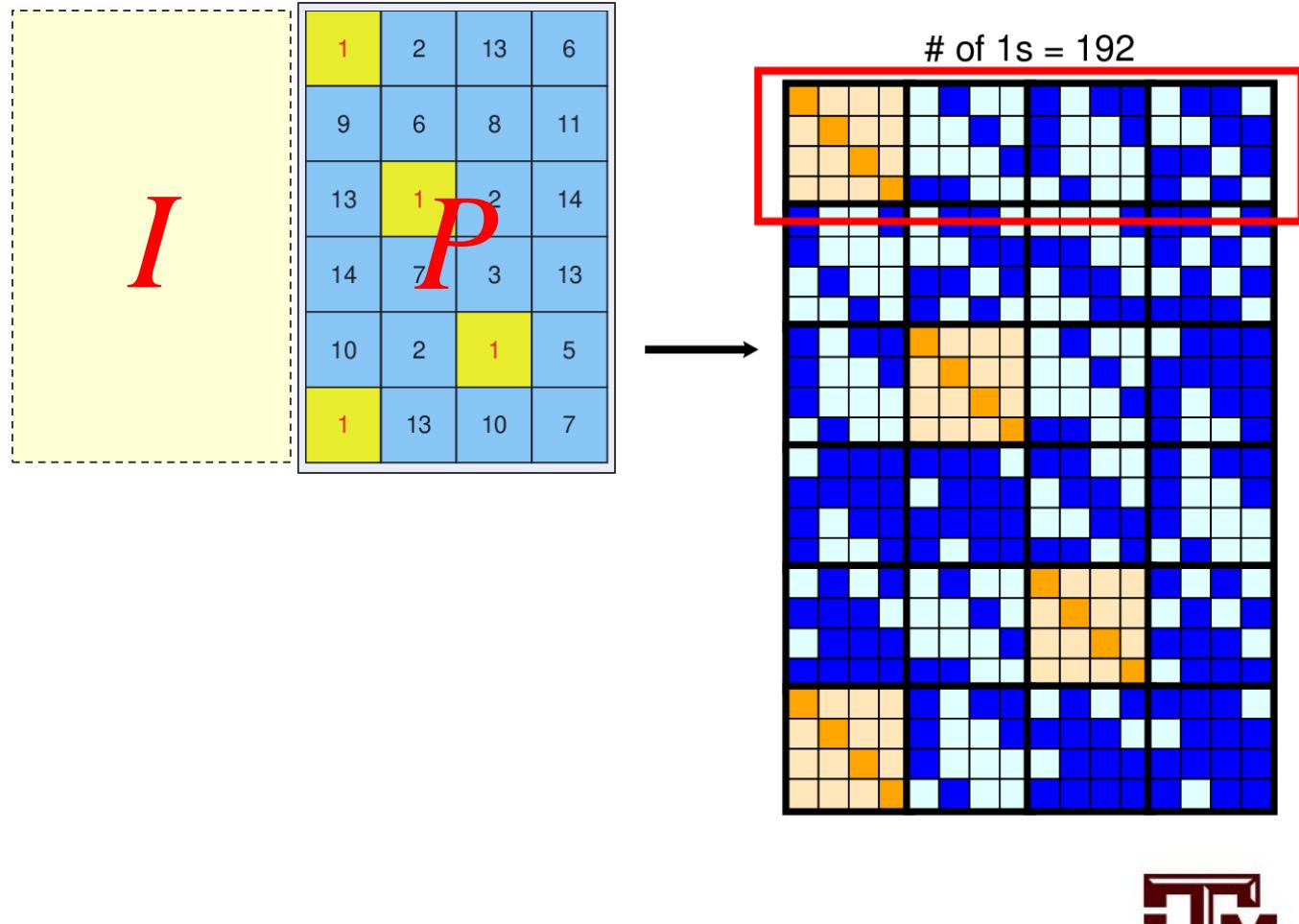
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



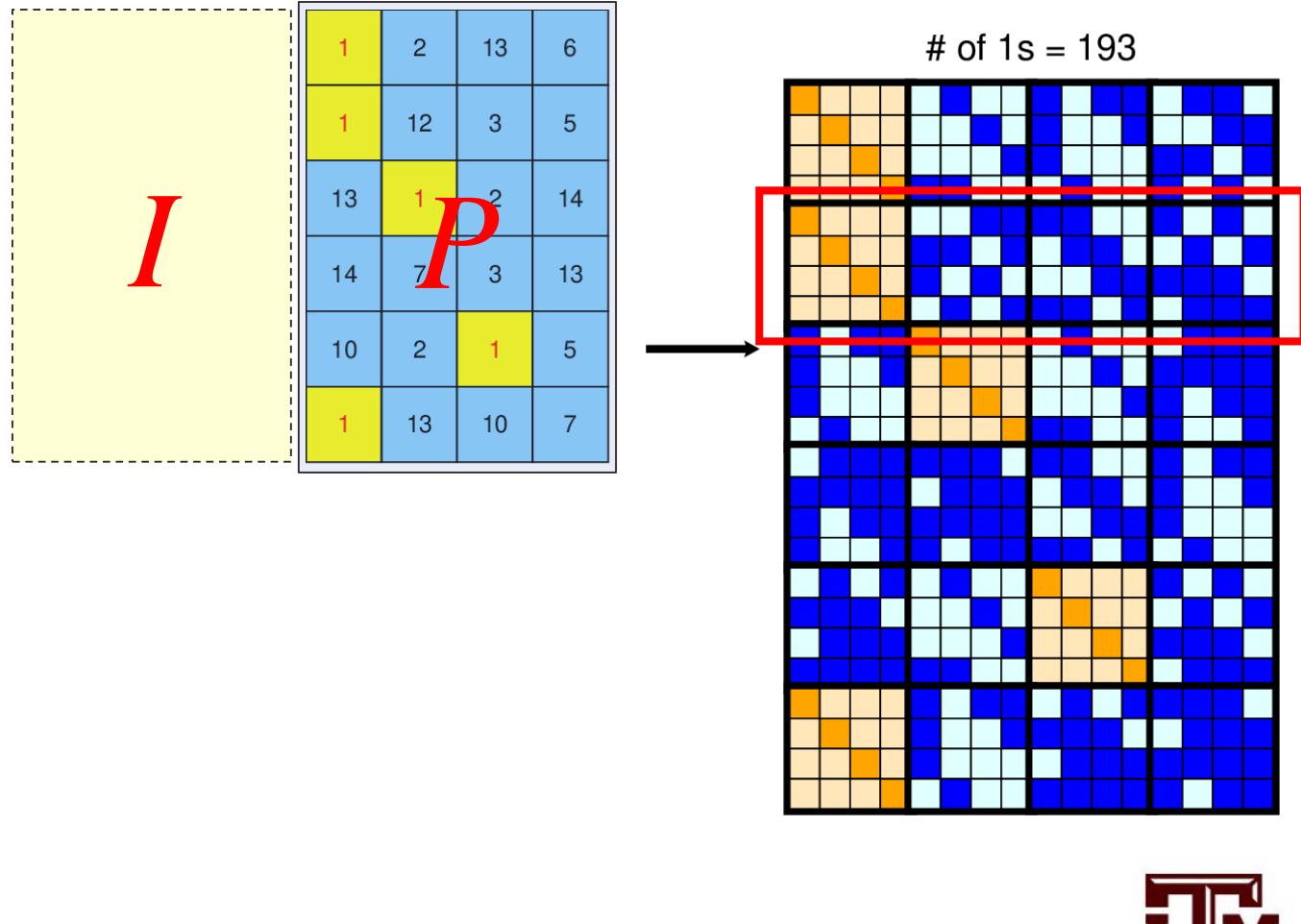
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



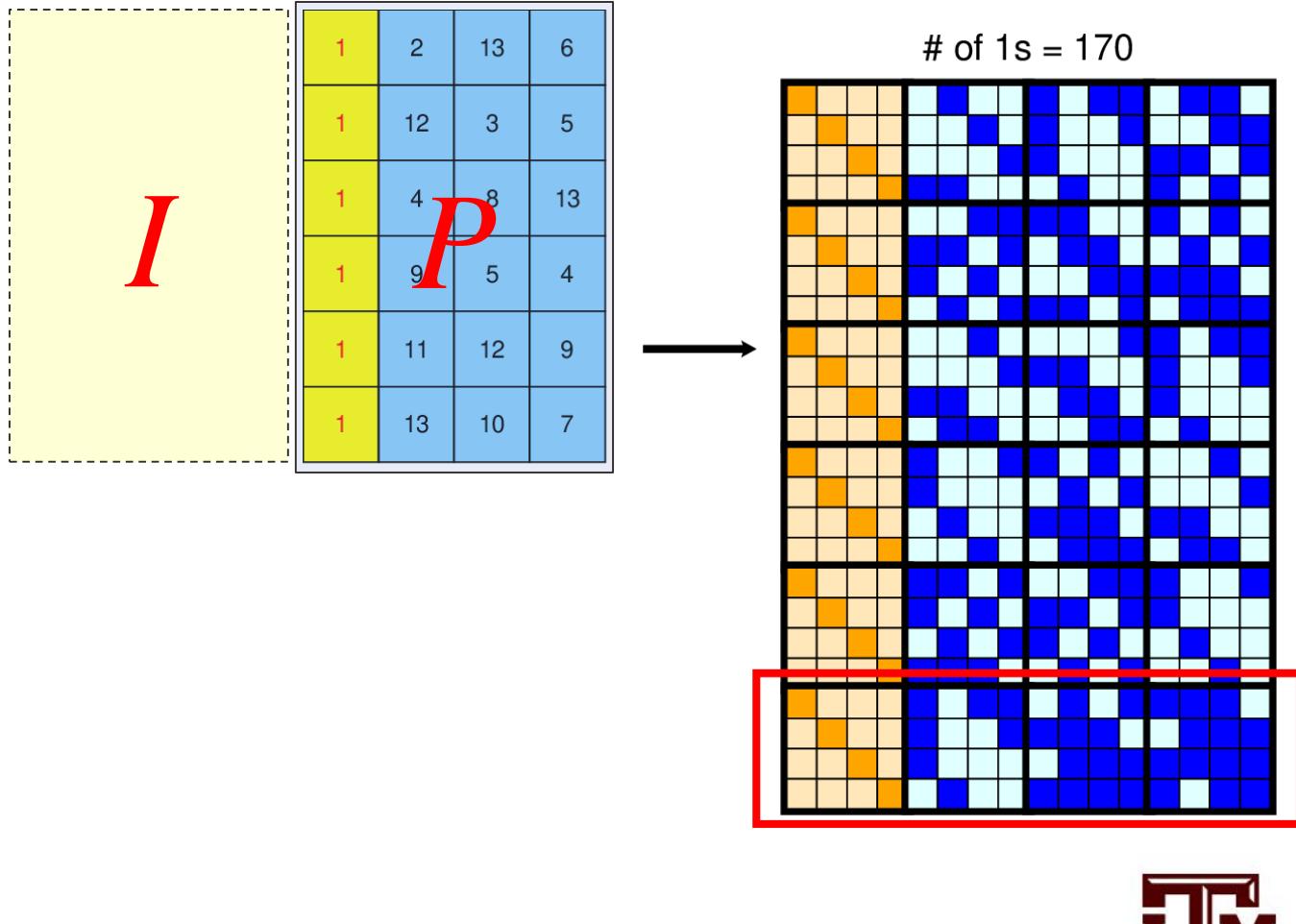
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



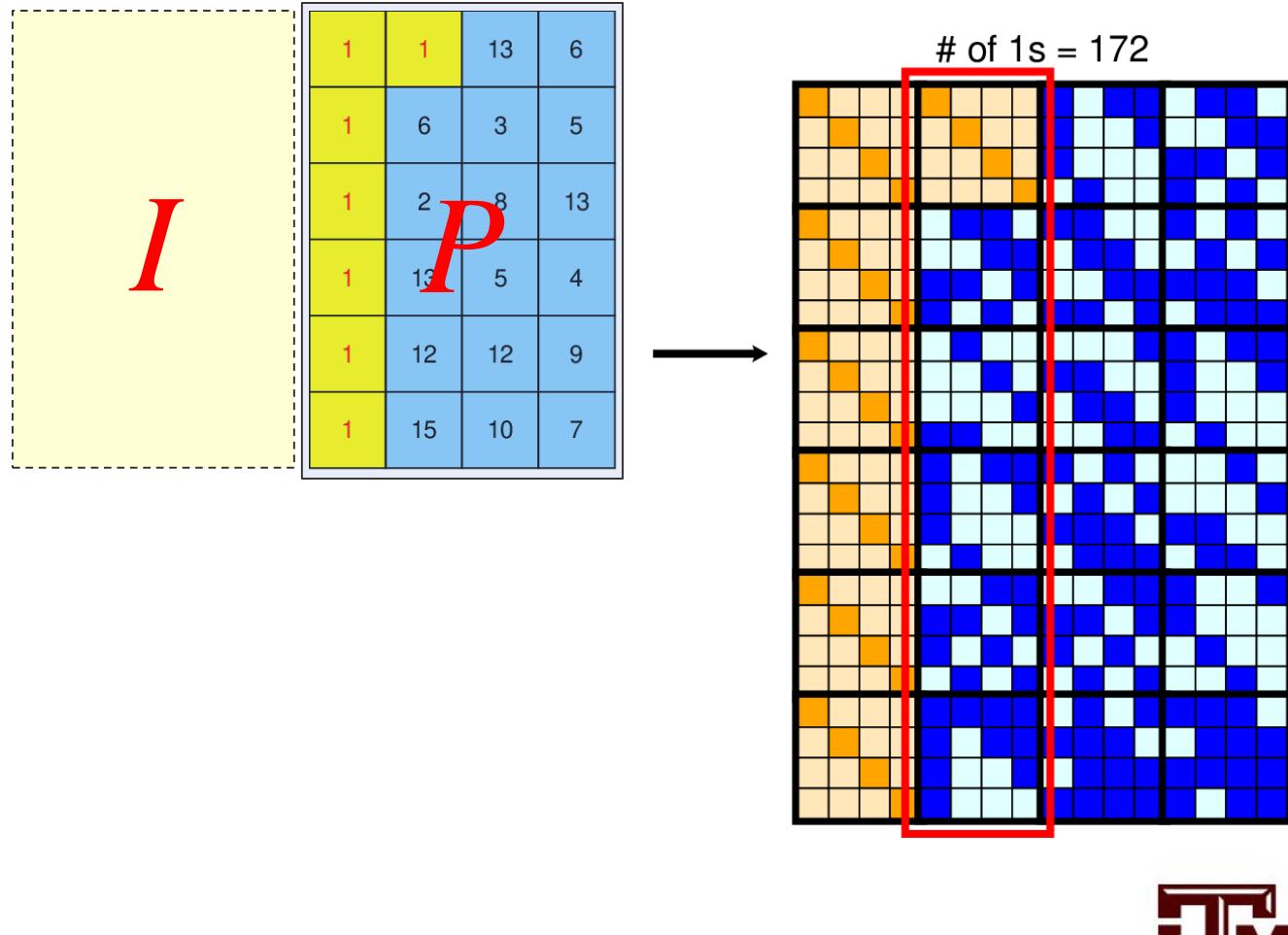
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



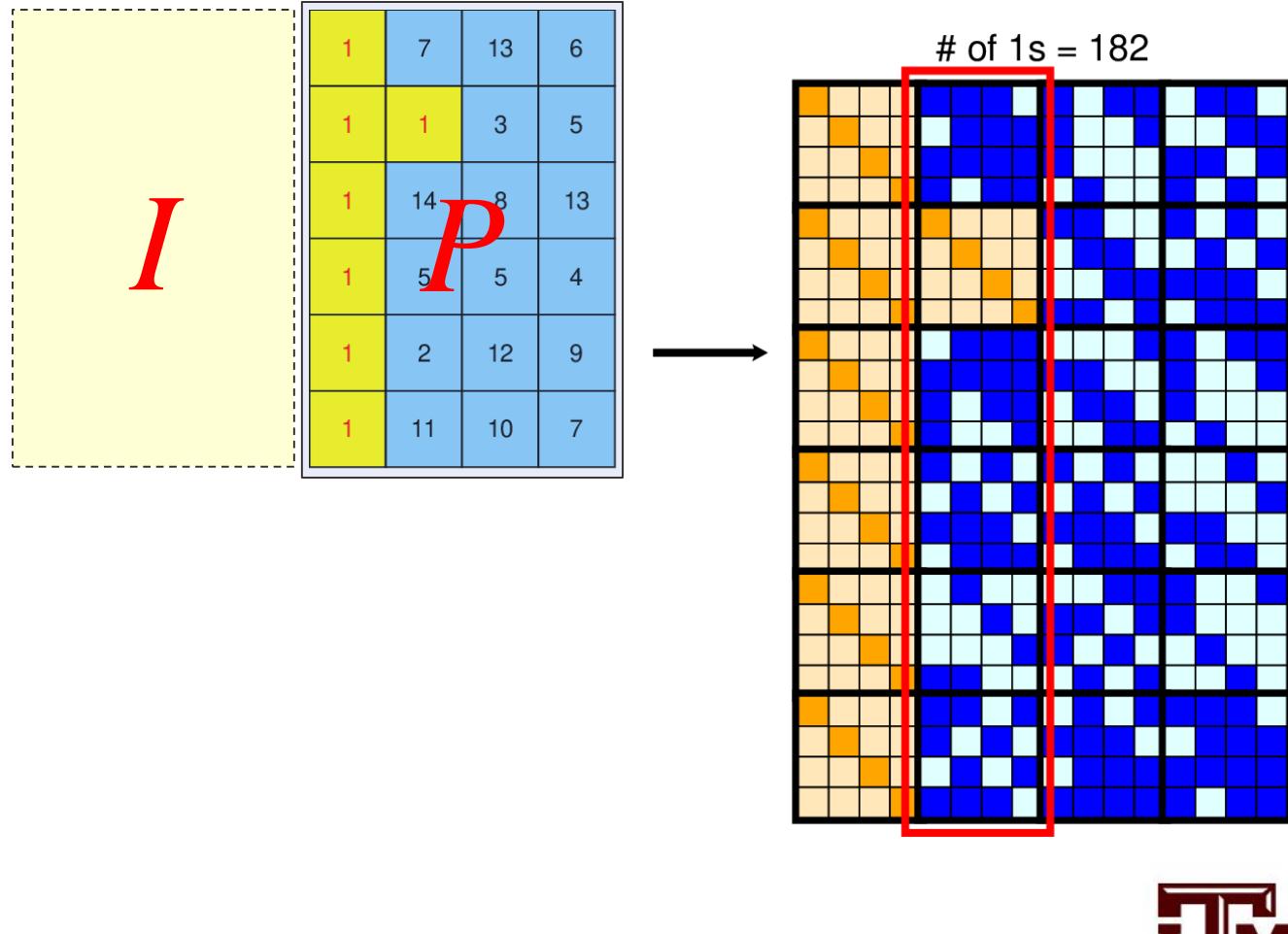
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



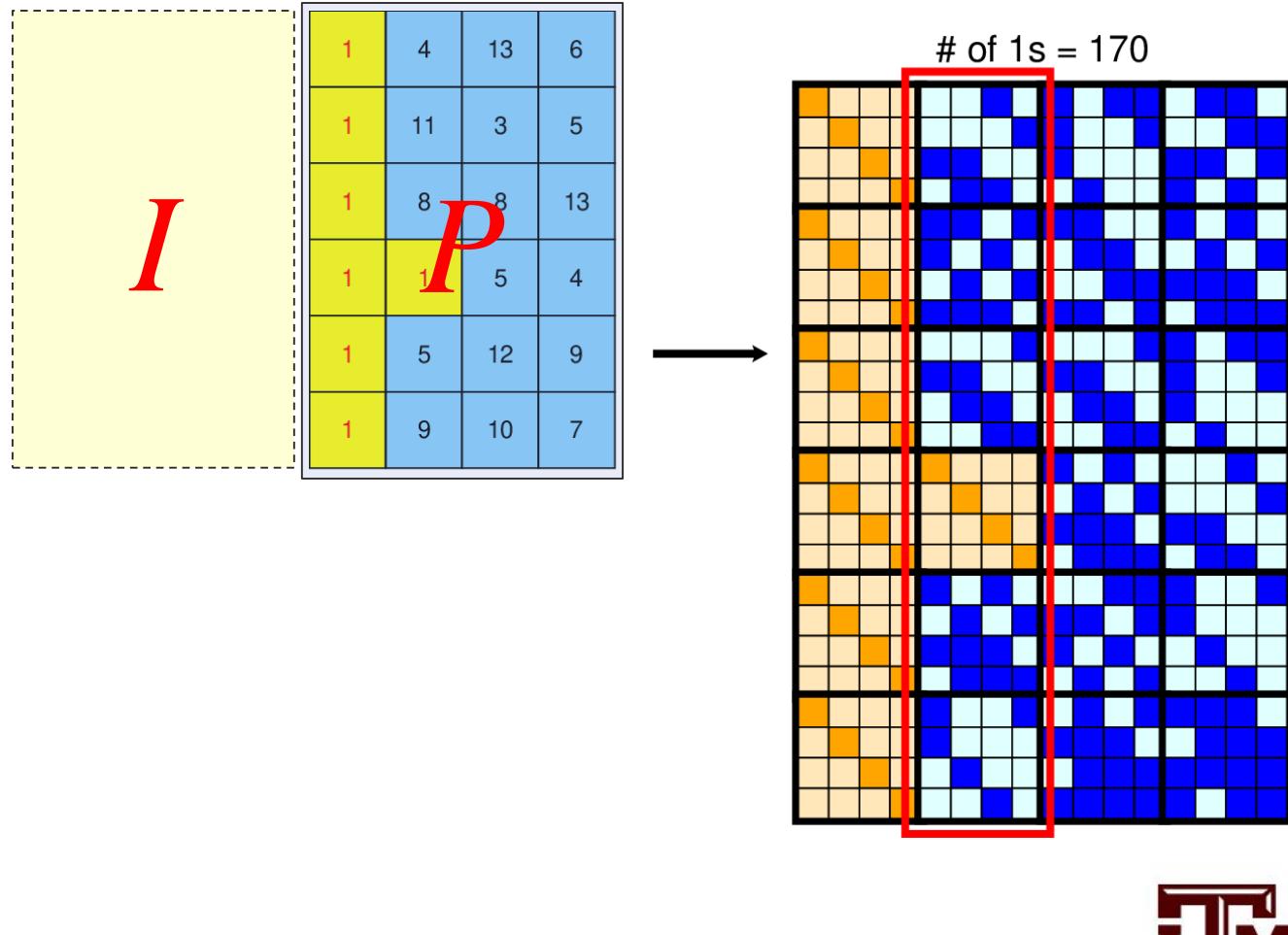
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



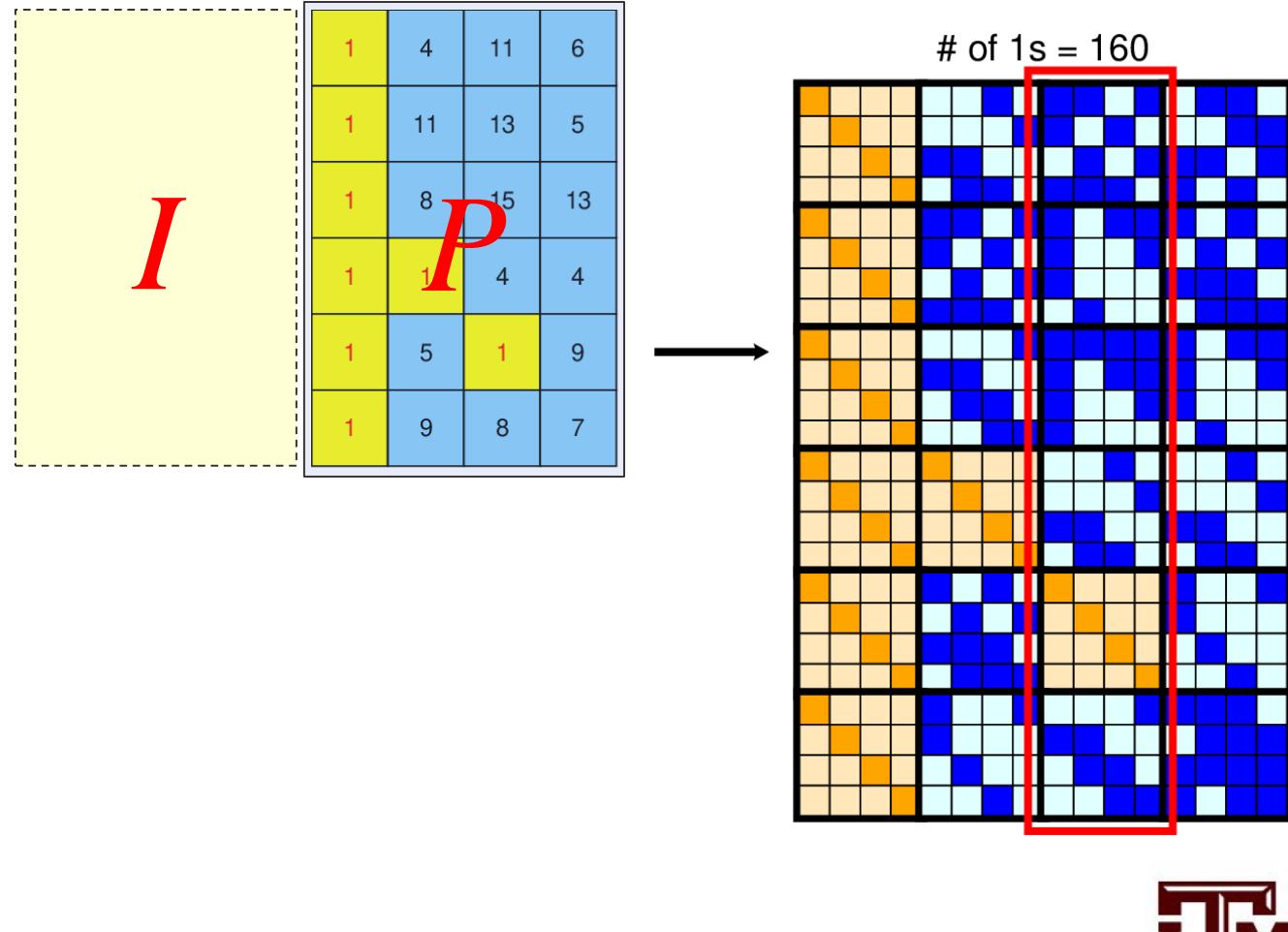
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



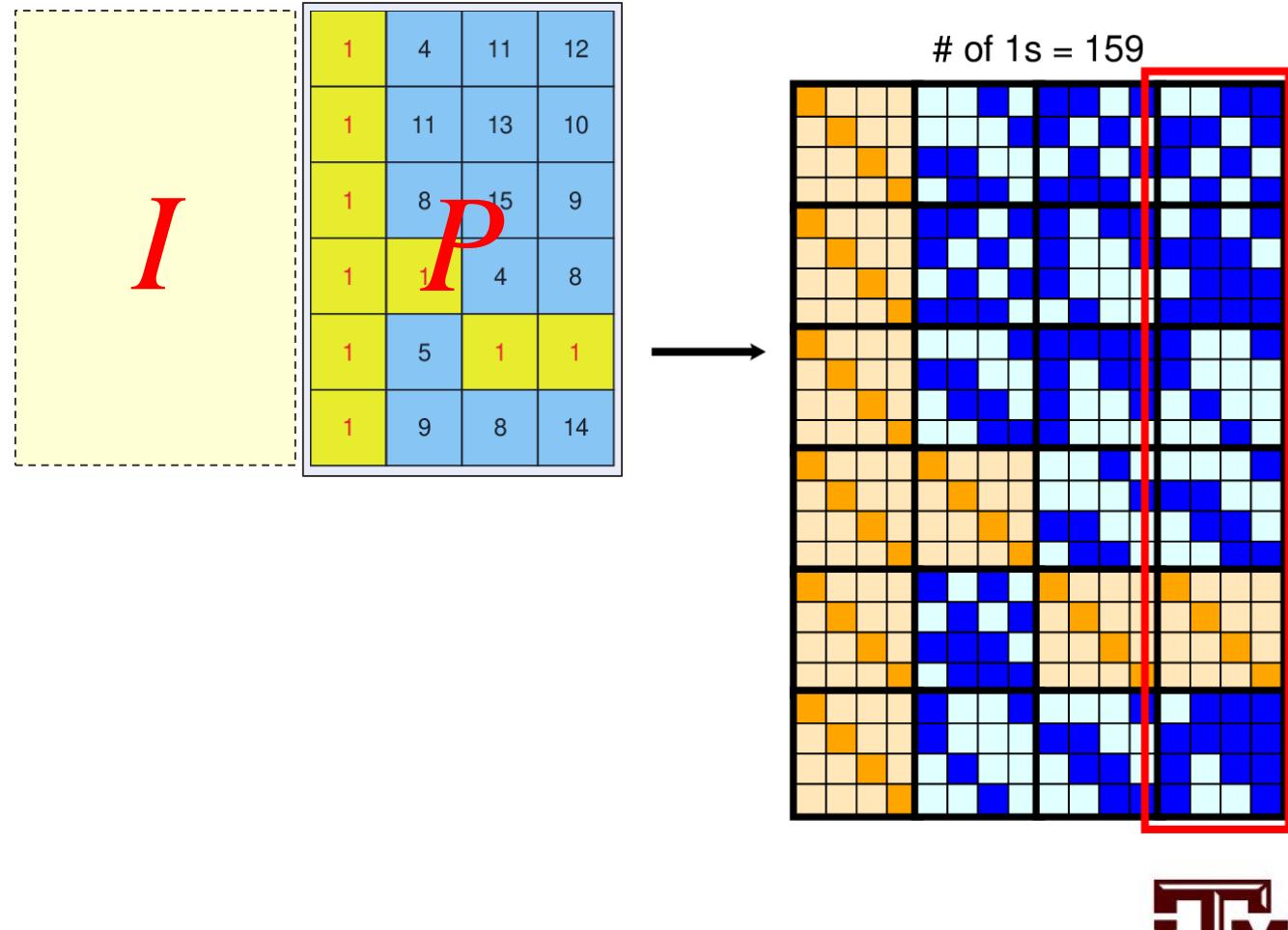
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



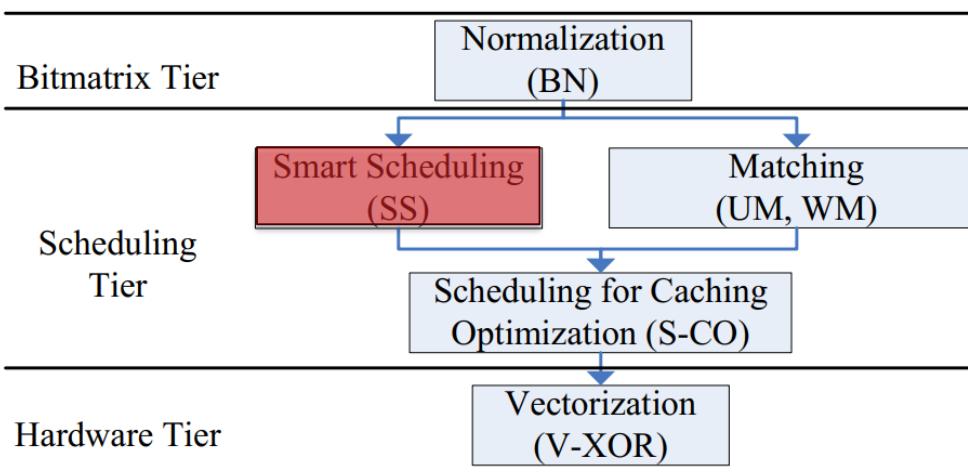
# Bitmatrix Normalization (BN)

- Normalize by each element, find less bit 1s



# Smart Scheduling (SS)

- Utilize both message and computed parities



# Smart Scheduling (SS)

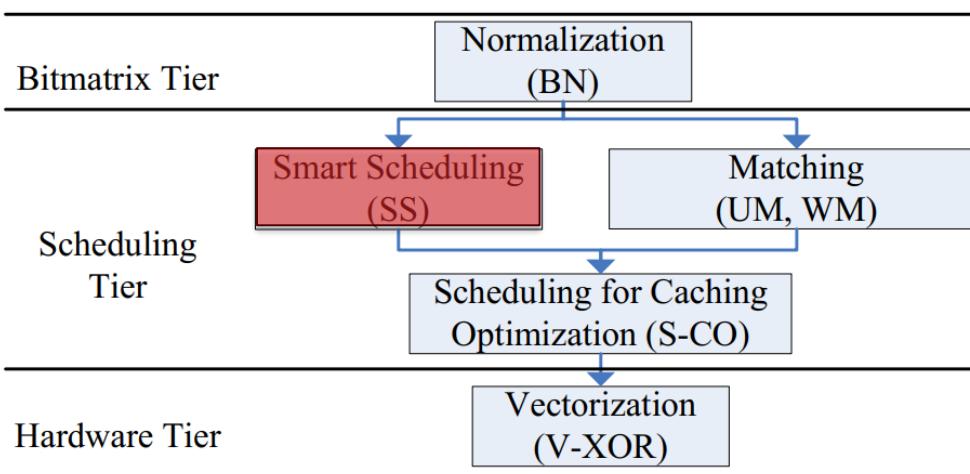
- Utilize both message and computed parities

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$


3 XORs

$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$


3 XORs



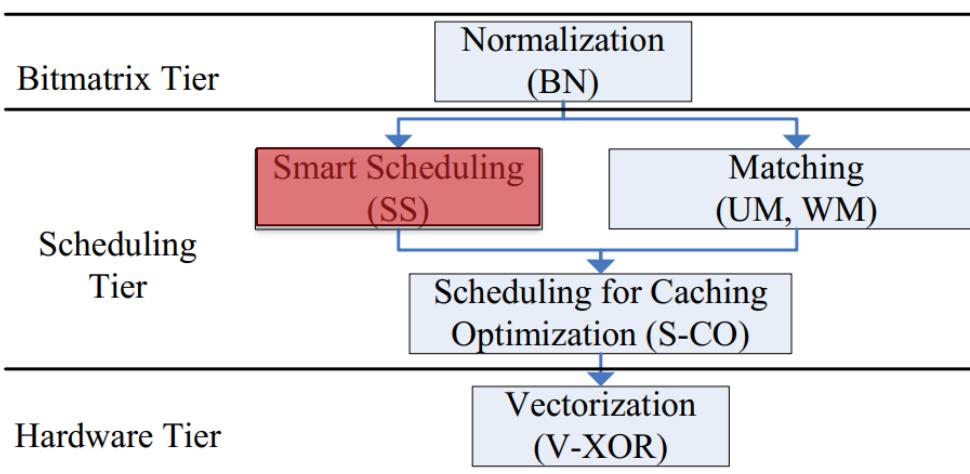
# Smart Scheduling (SS)

- Utilize both message and computed parities

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

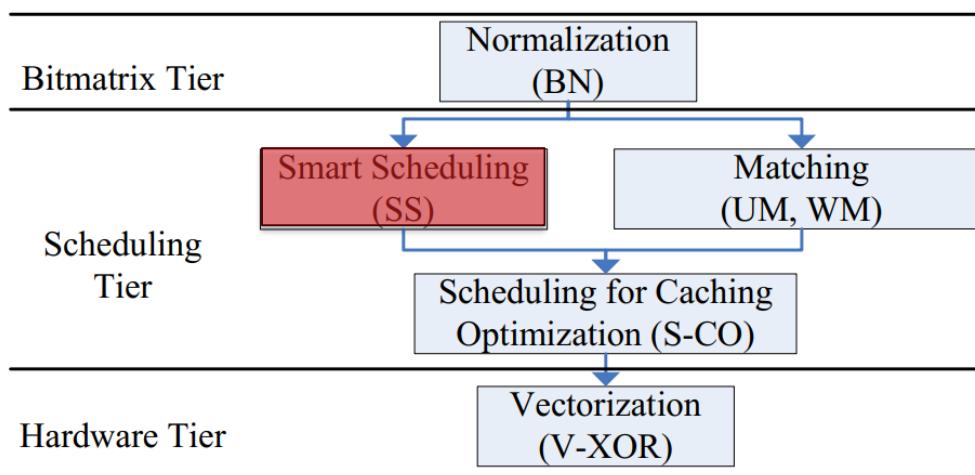

3 XORs

3 XORs

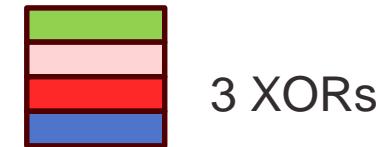


# Smart Scheduling (SS)

- Utilize both message and computed parities



$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

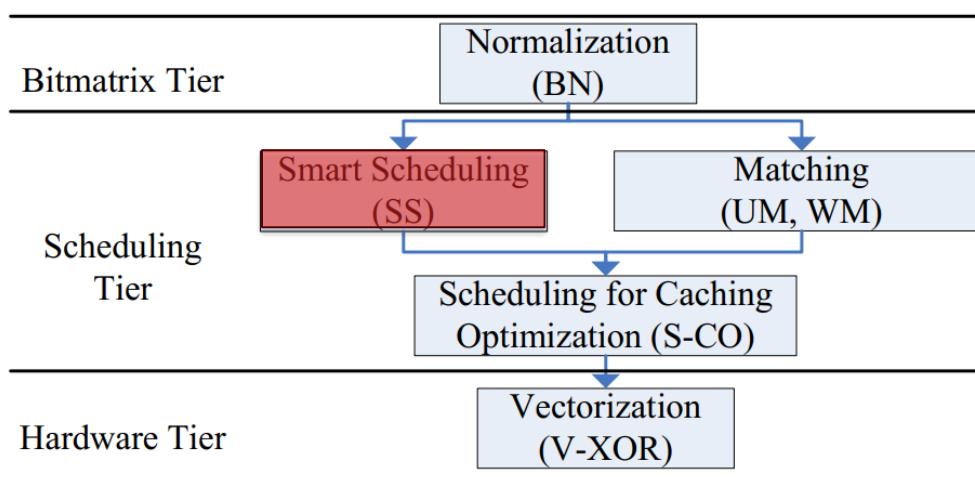


$$p_1 = p_0 \oplus u_4 \oplus u_5$$



# Smart Scheduling (SS)

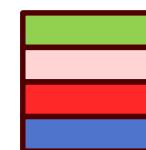
- Utilize both message and computed parities



$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

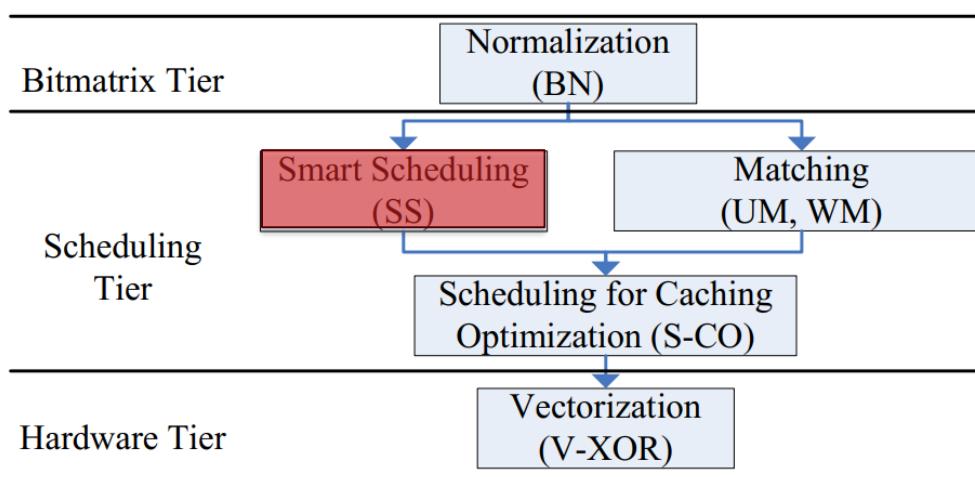


$$p_1 = p_0 \oplus u_4 \oplus u_5$$



# Smart Scheduling (SS)

- Utilize both message and computed parities



$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

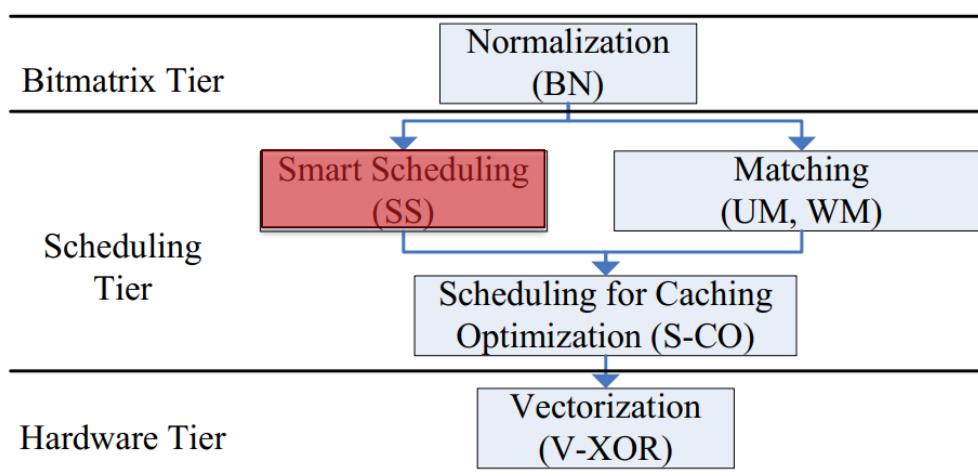


$$p_1 = p_0 \oplus u_4 \oplus u_5$$

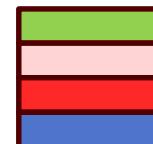


# Smart Scheduling (SS)

- Utilize both message and computed parities



$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$



3 XORs

$$p_1 = p_0 \oplus u_4 \oplus u_5$$



3 XORs

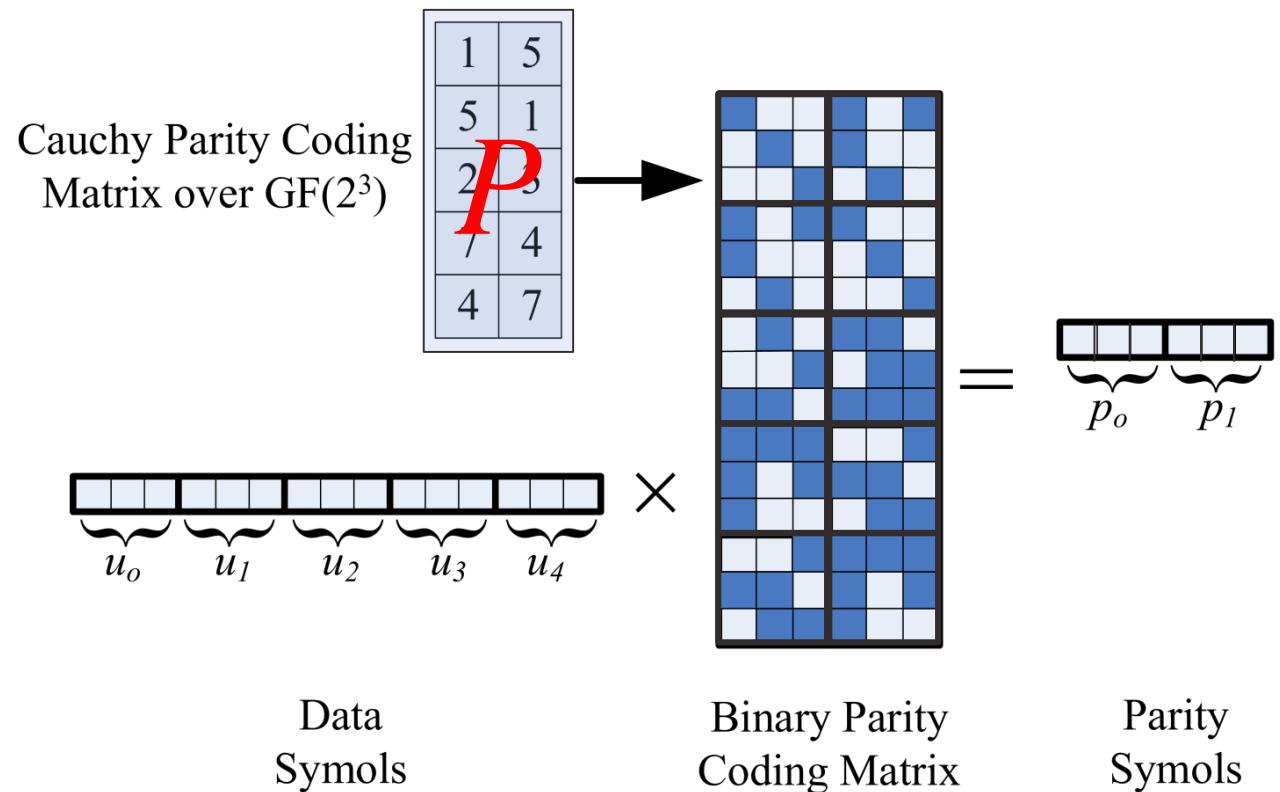
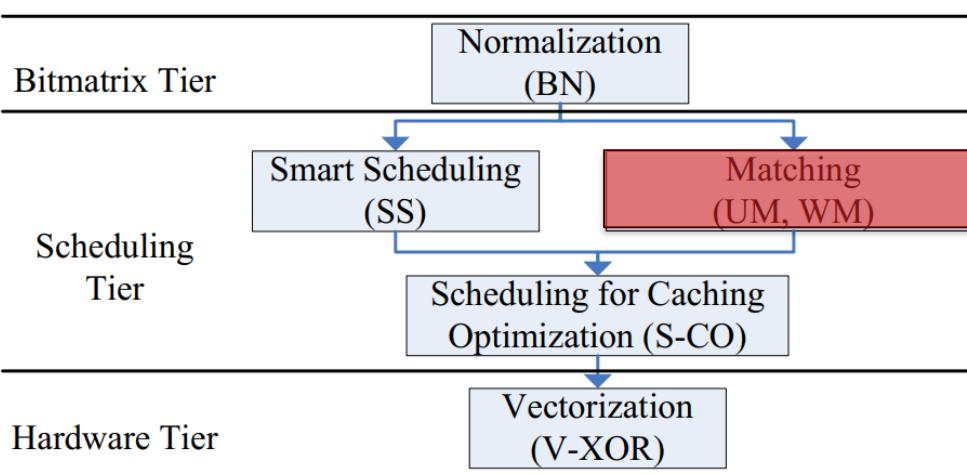


2 XORs



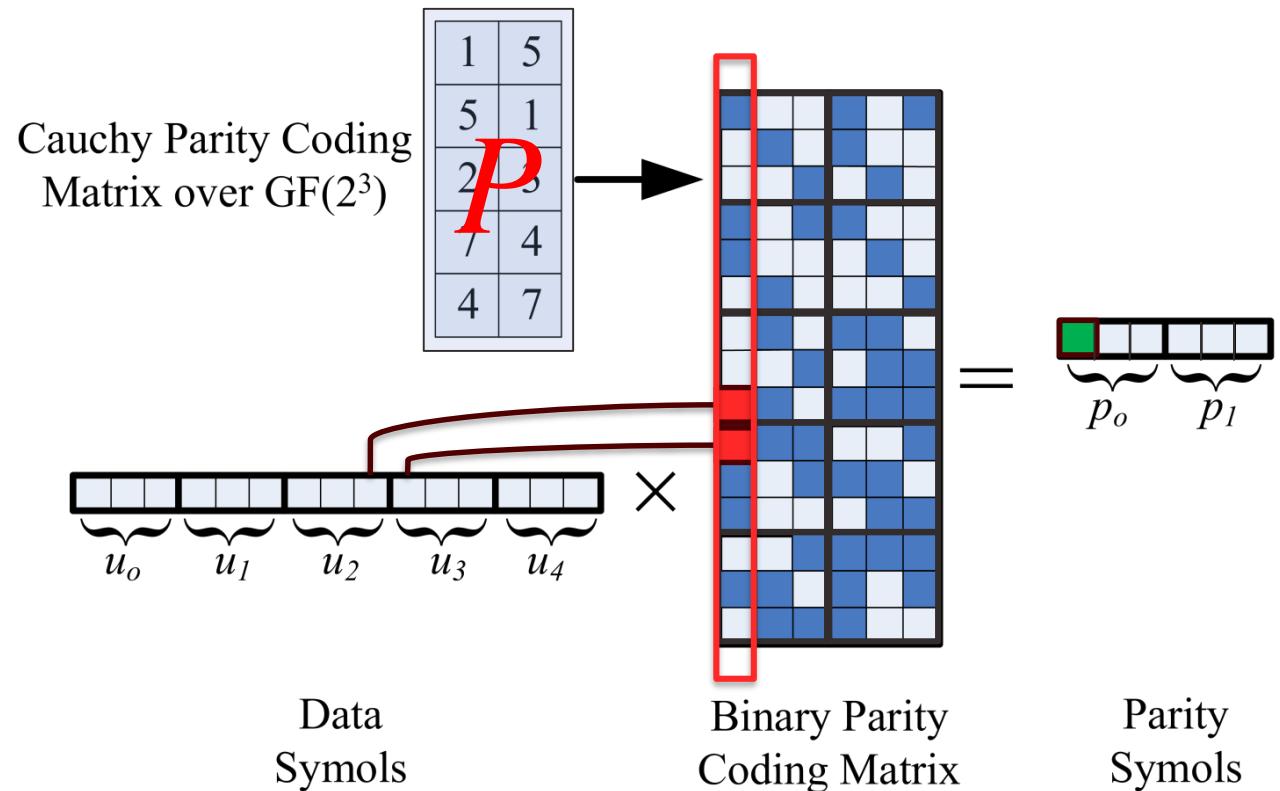
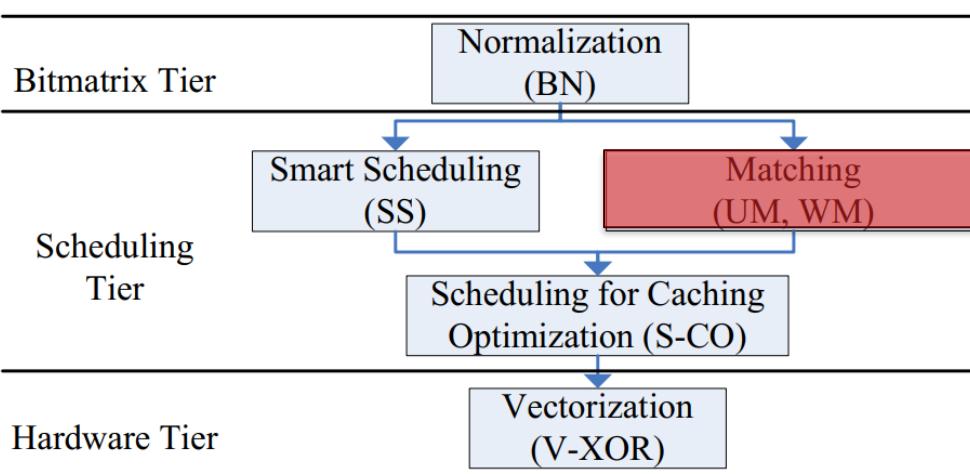
# Matching (UM, WM)

- Common XORs of pair message bits
- Unweighted/weighted greedy algorithm



# Matching (UM, WM)

- Common XORs of pair message bits
- Unweighted/weighted greedy algorithm

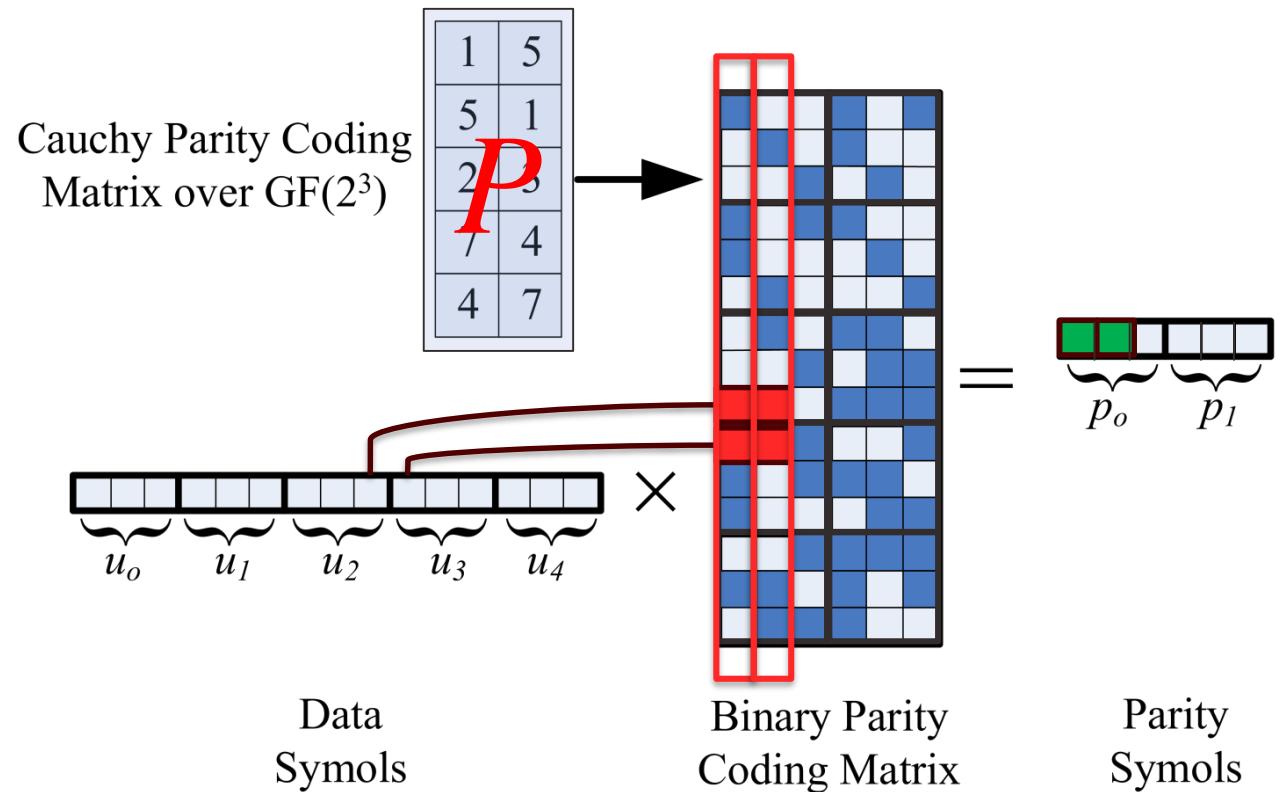
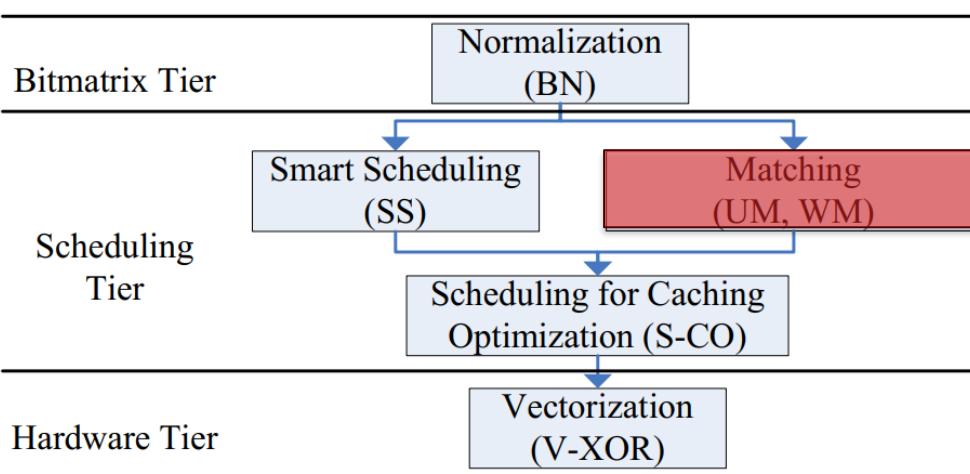


Cheng Huang, Jin Li, and Minghua Chen. On optimizing XOR-based codes for fault-tolerant storage applications. In Proceedings of Information Theory Workshop, pages 218–223, 2007.



# Matching (UM, WM)

- Common XORs of pair message bits
- Unweighted/weighted greedy algorithm

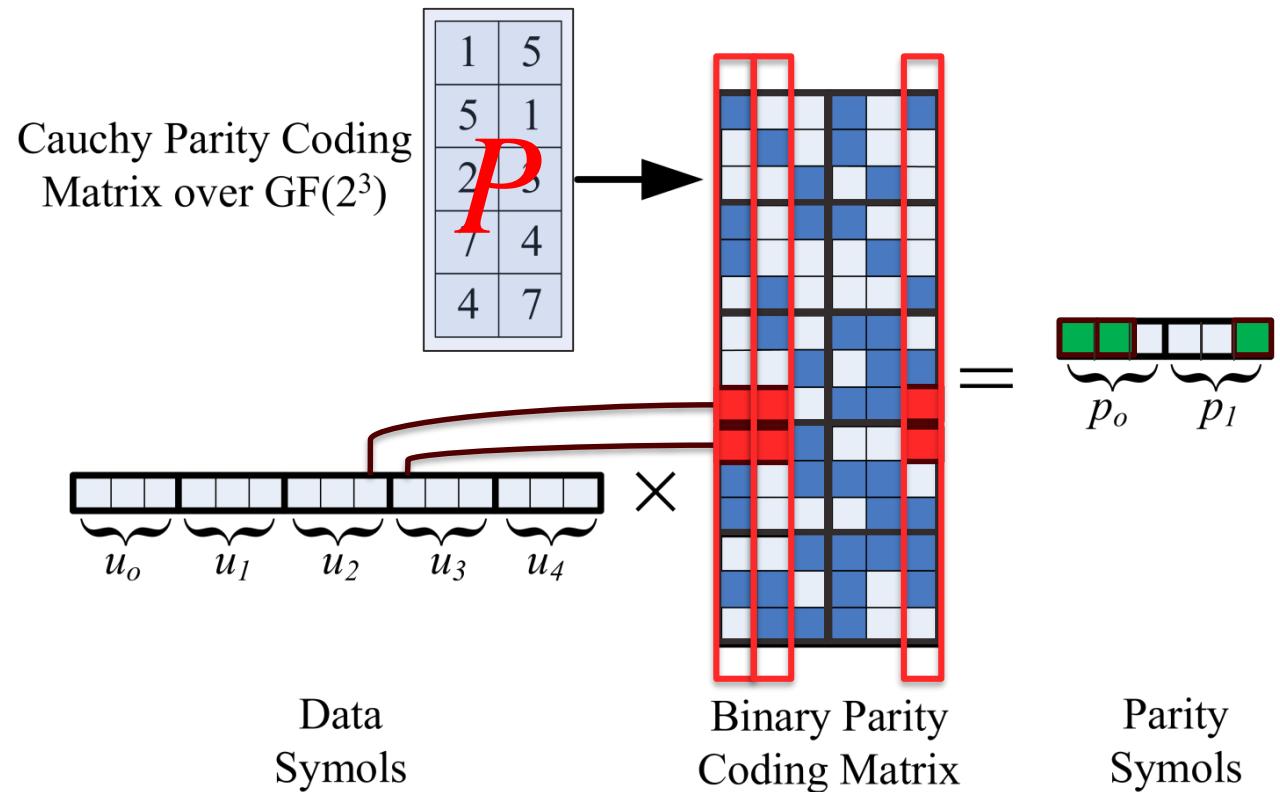
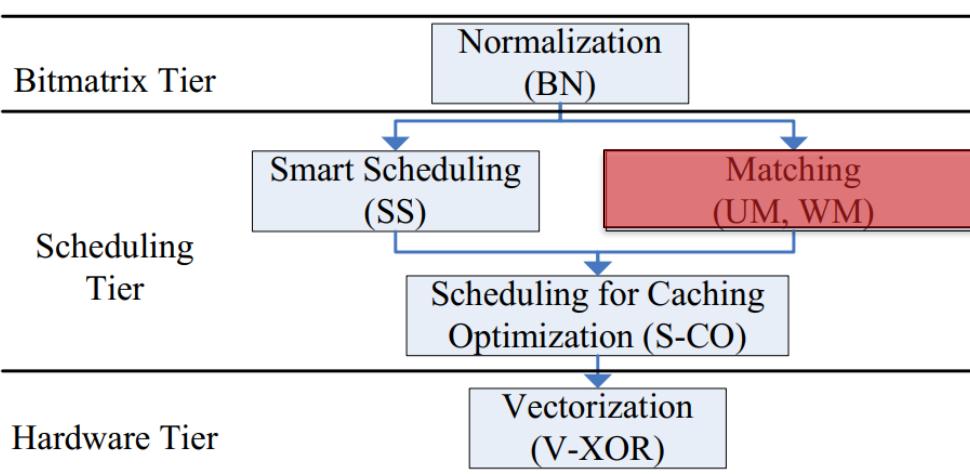


Cheng Huang, Jin Li, and Minghua Chen. On optimizing XOR-based codes for fault-tolerant storage applications. In Proceedings of Information Theory Workshop, pages 218–223, 2007.



# Matching (UM, WM)

- Common XORs of pair message bits
- Unweighted/weighted greedy algorithm

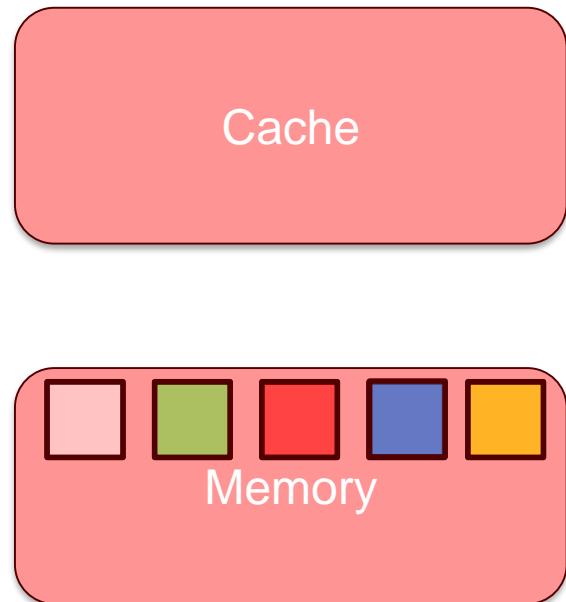
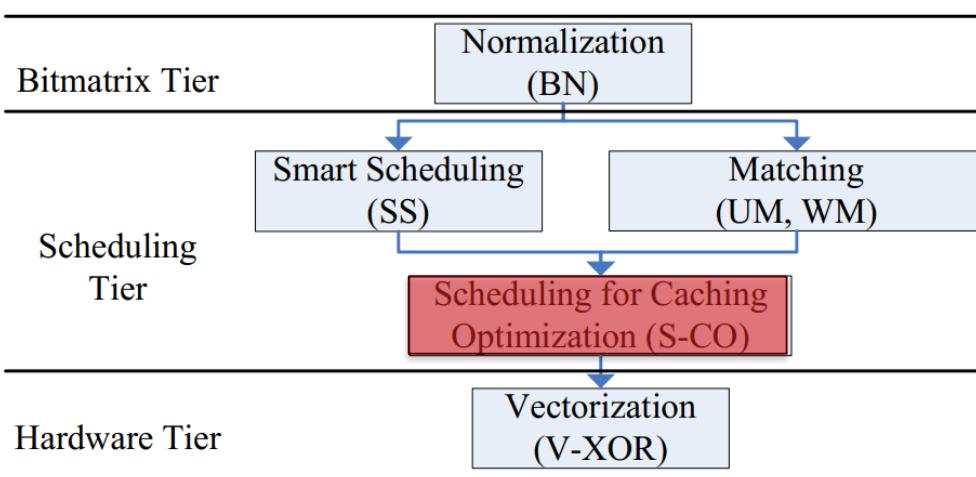


# Scheduling - Cache Optimization (S-CO)

- Reduce cache miss penalty

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$

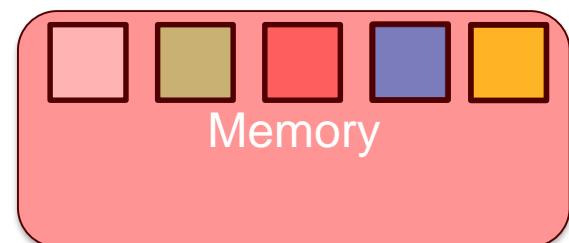
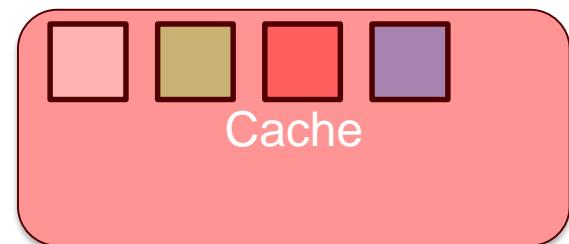
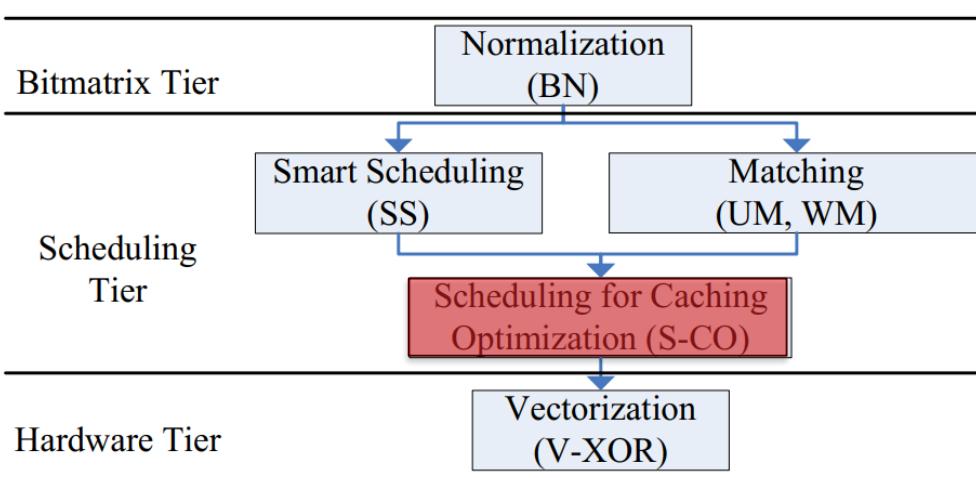


# Scheduling - Cache Optimization (S-CO)

- Reduce cache miss penalty

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$

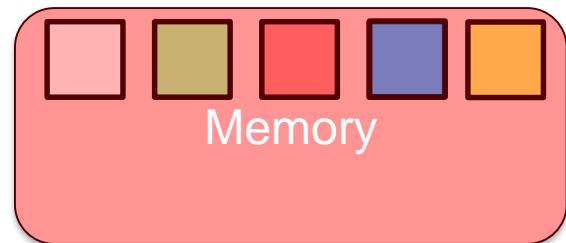
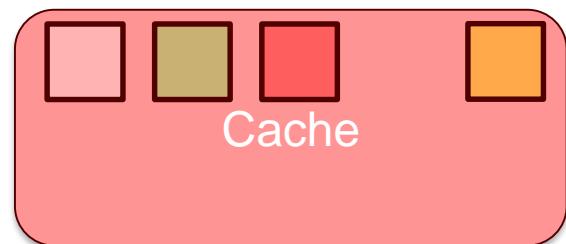
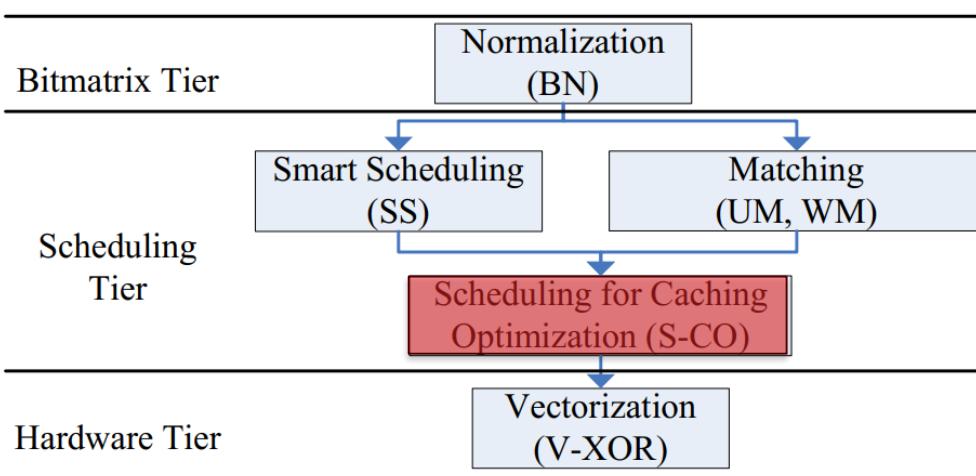


# Scheduling - Cache Optimization (S-CO)

- Reduce cache miss penalty

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

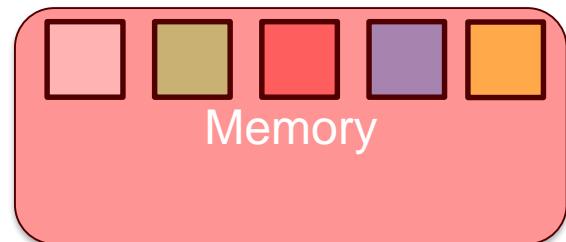
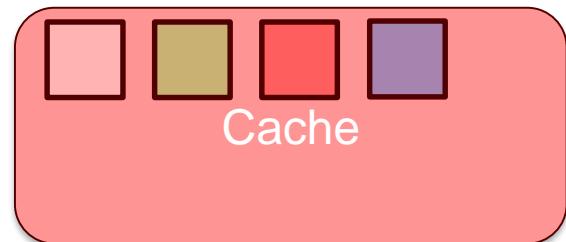
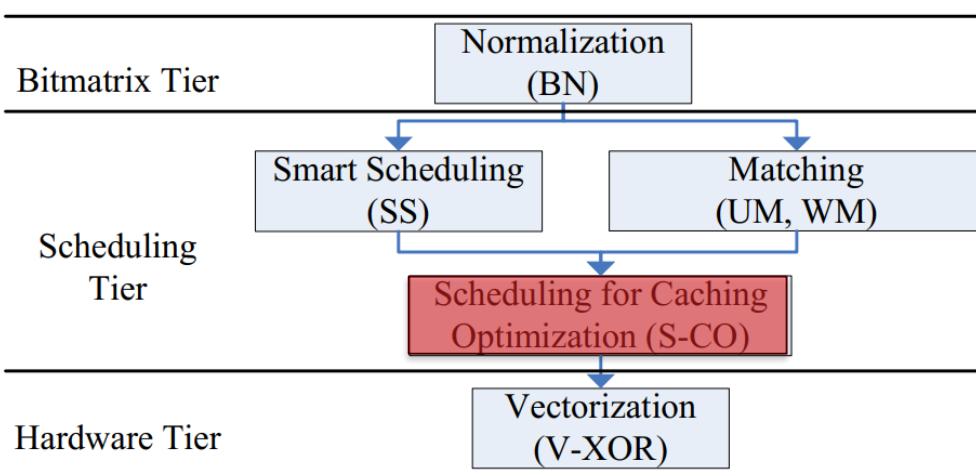
$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$



# Scheduling - Cache Optimization (S-CO)

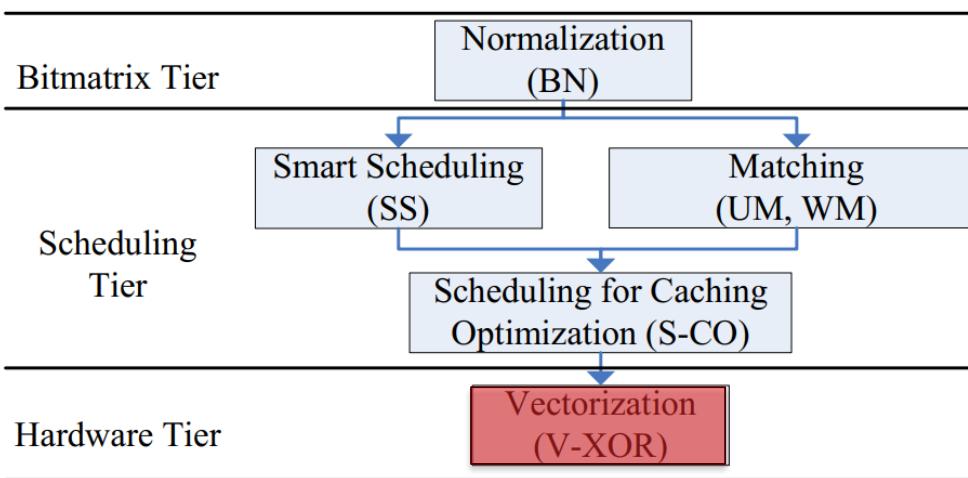
- Reduce cache miss penalty

$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$
$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$



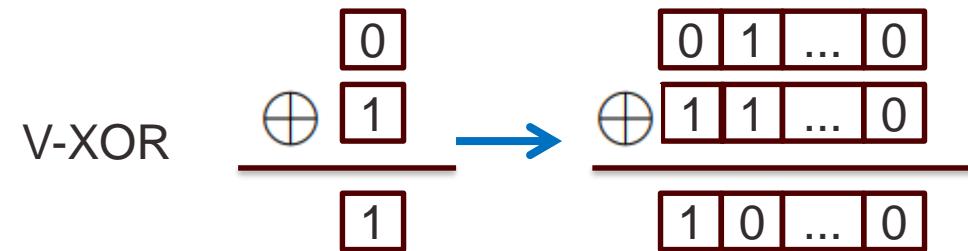
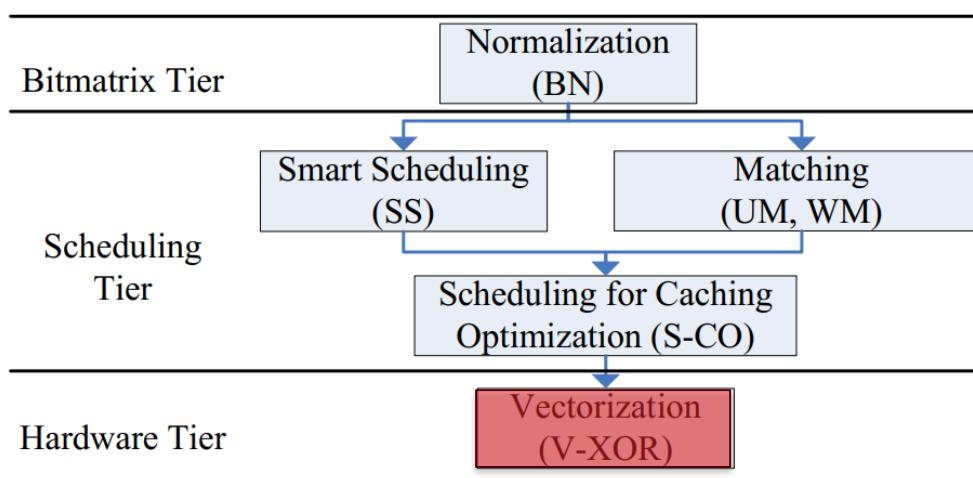
# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits



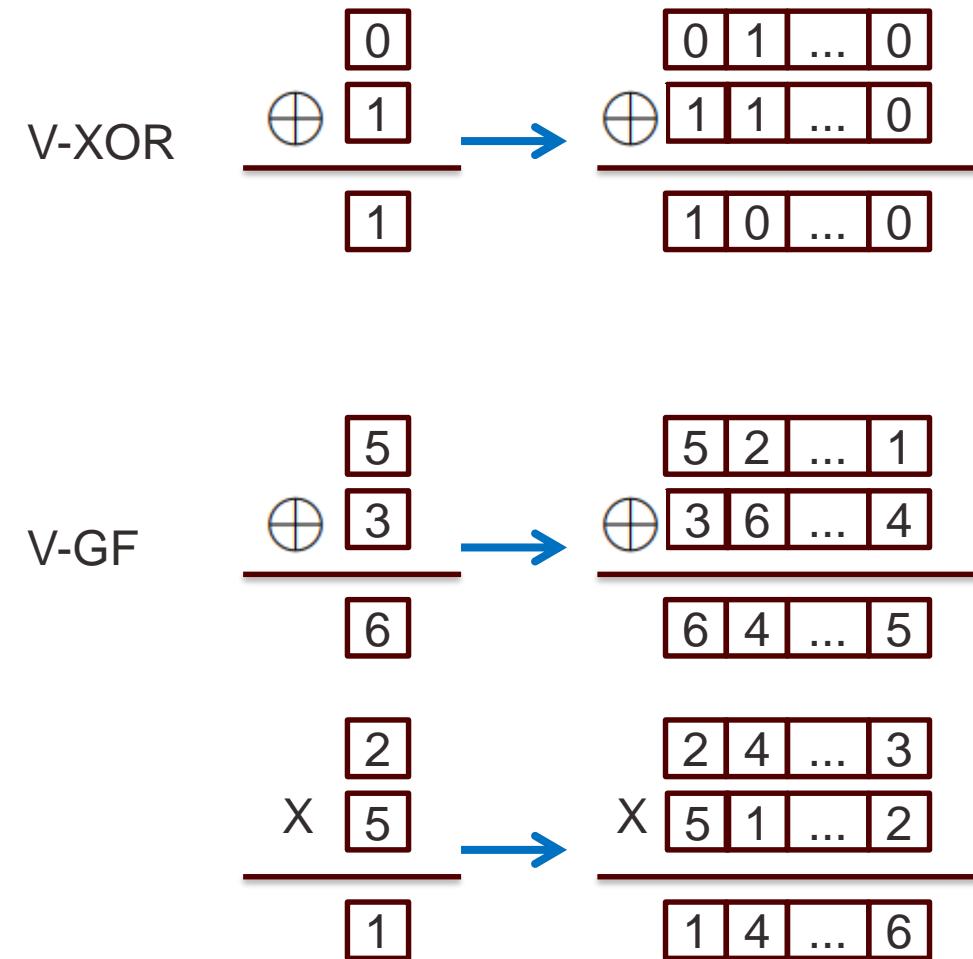
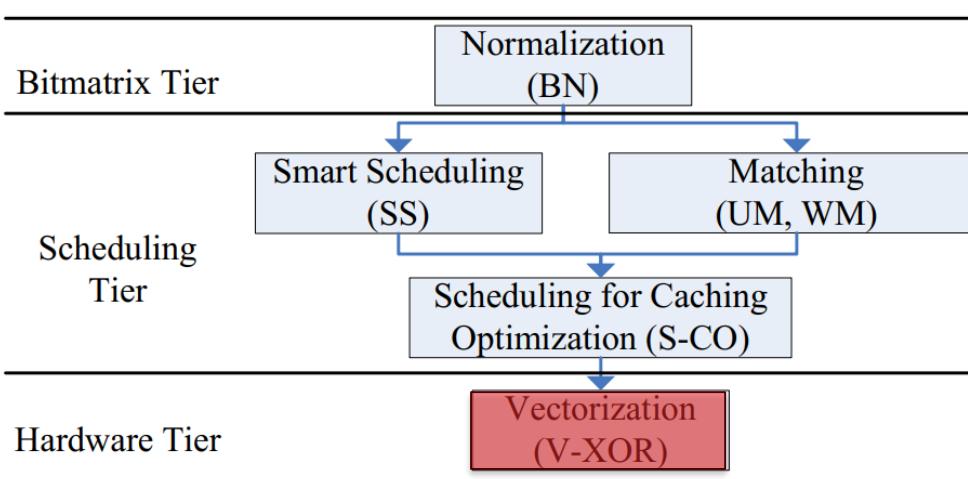
# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits



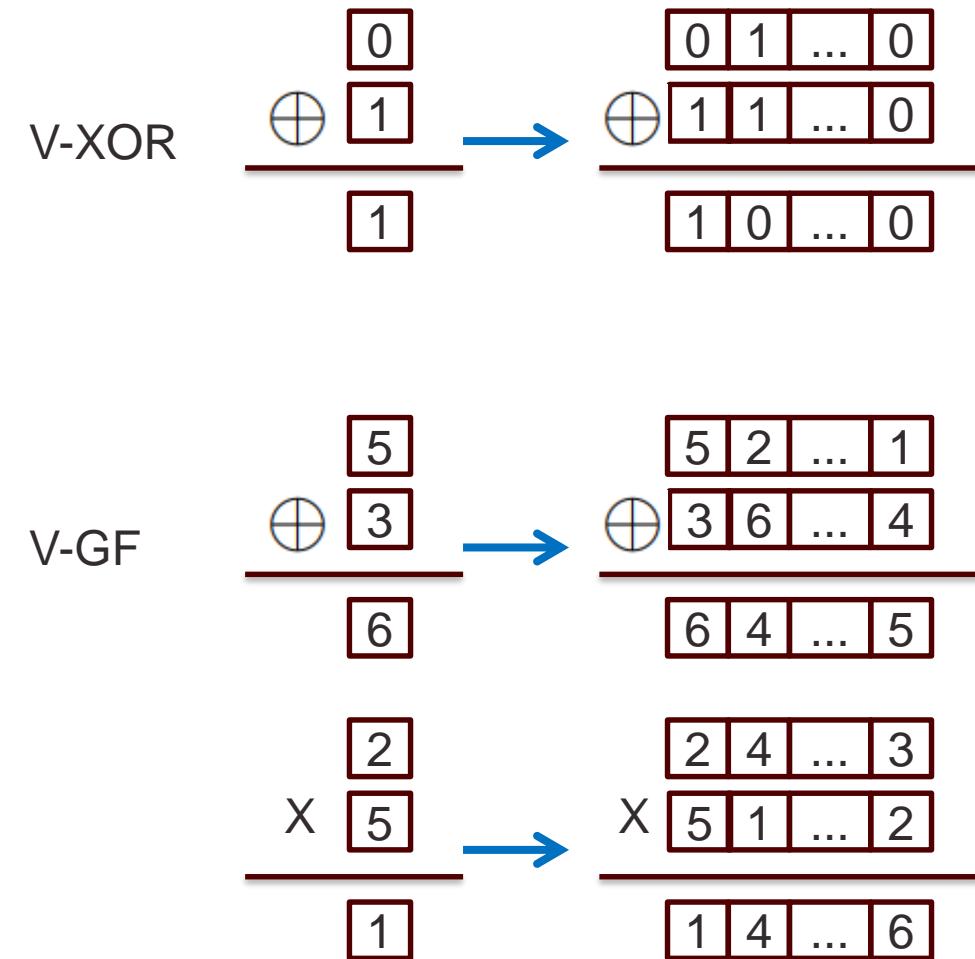
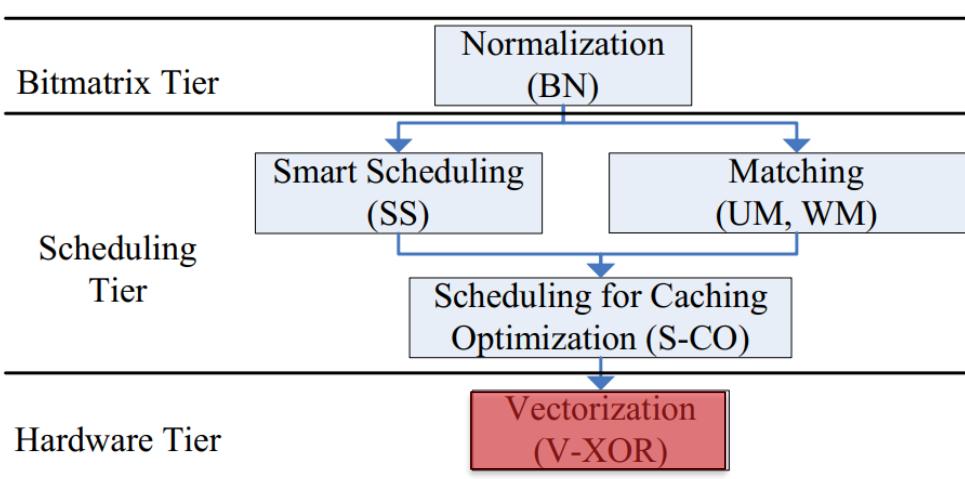
# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits



# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits

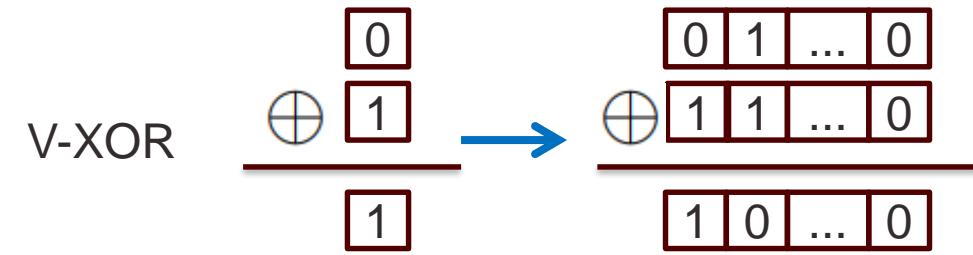
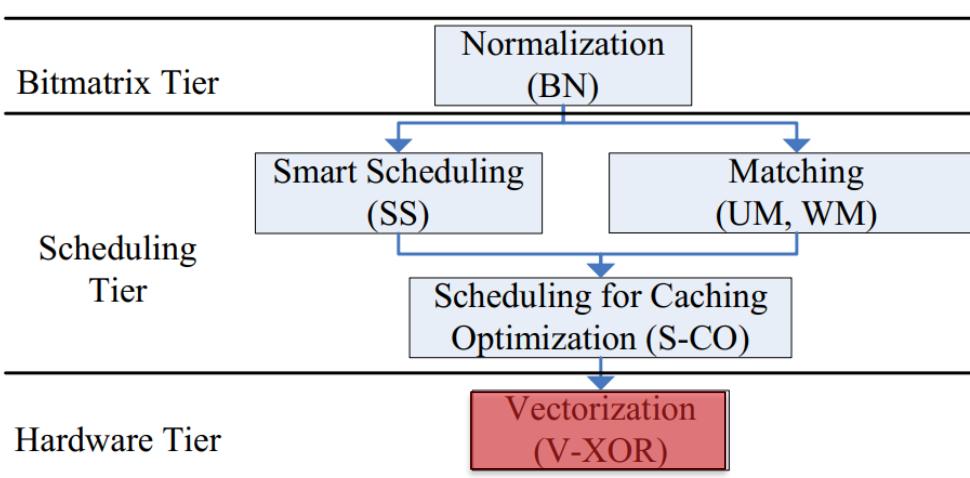


**Jerasure 2.0**

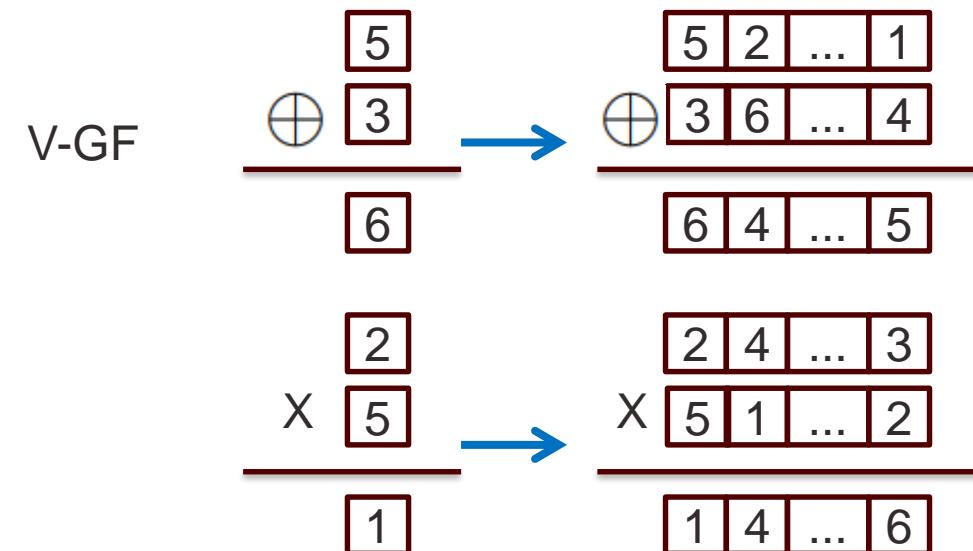


# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits



**V.S.**

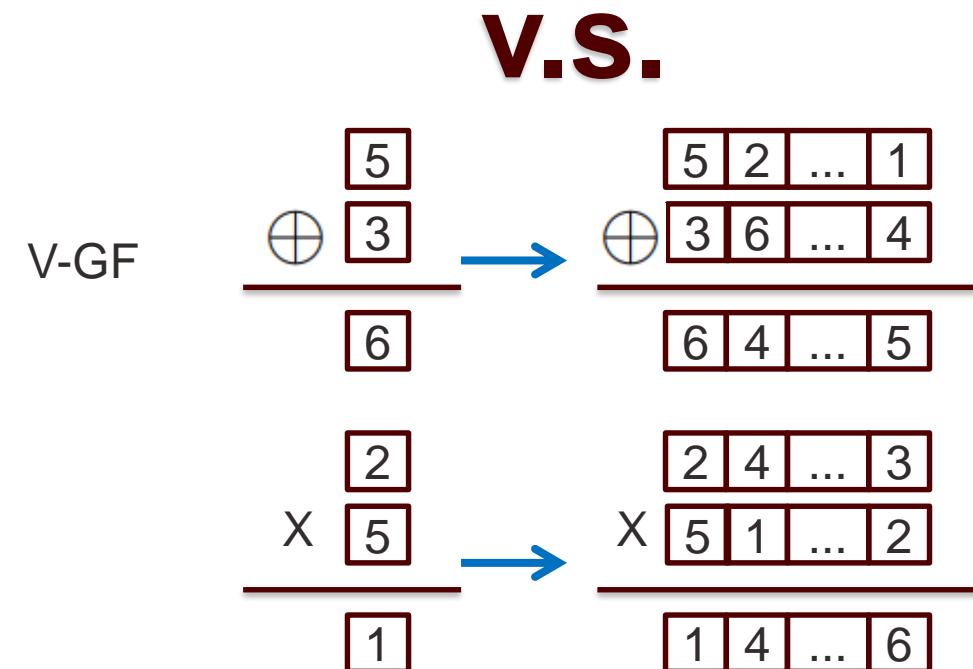
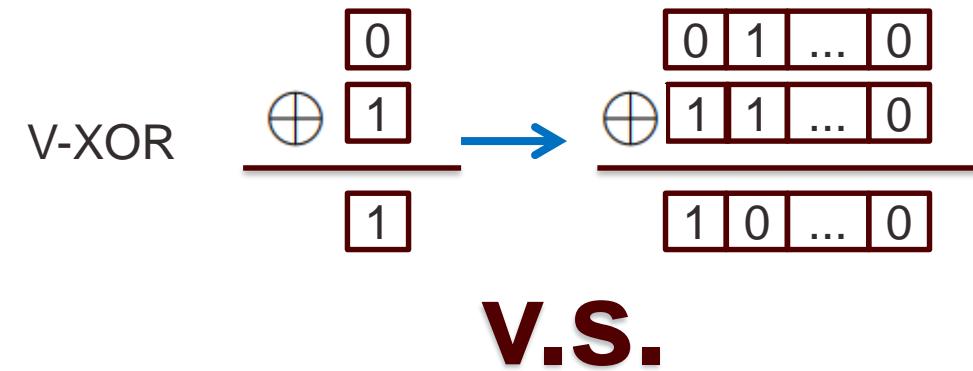
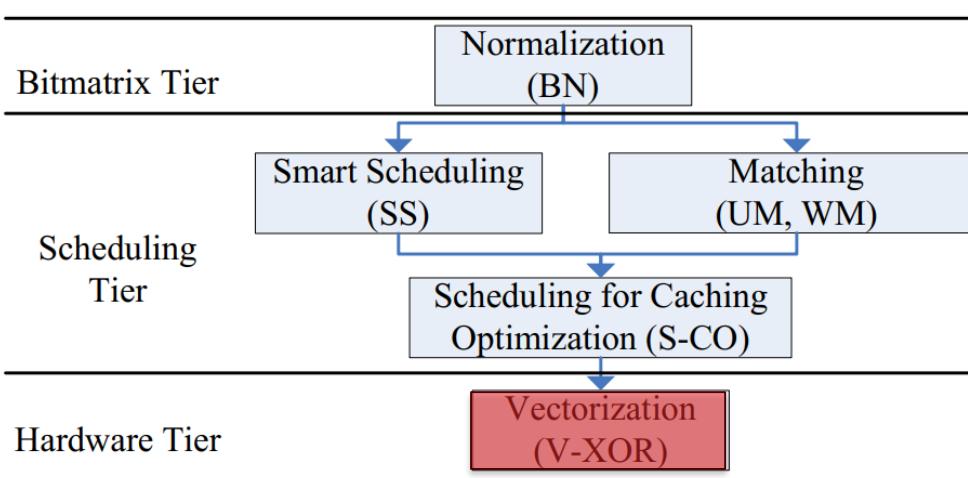


**Jerasure 2.0**



# Vectorization

- Using SIMD ISA perform multiple bits operation
  - SSE: 128 bits
  - AVX2: 256 bits
  - AVX512: 512 bits



**Jerasure 2.0**

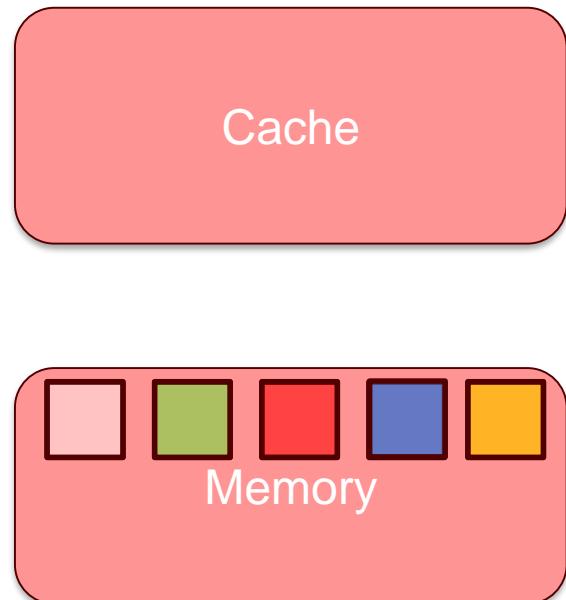
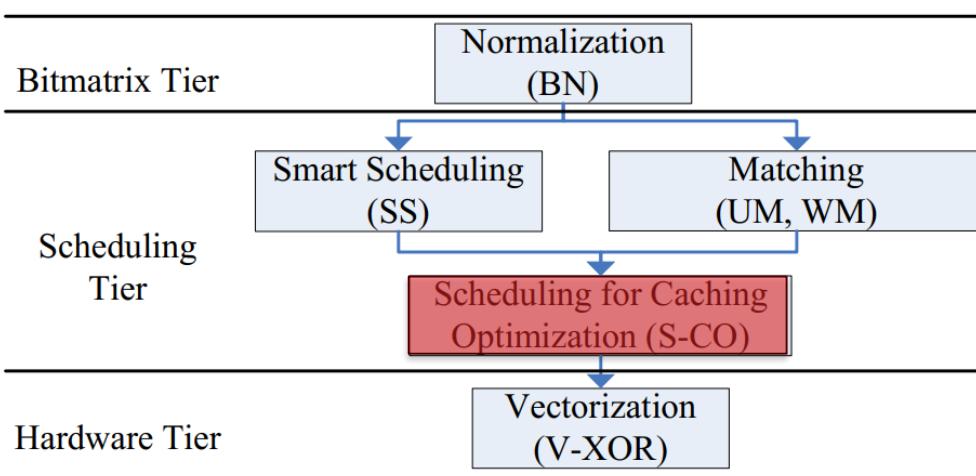


# Scheduling - Cache Optimization (S-CO)

- Reduce cache miss penalty

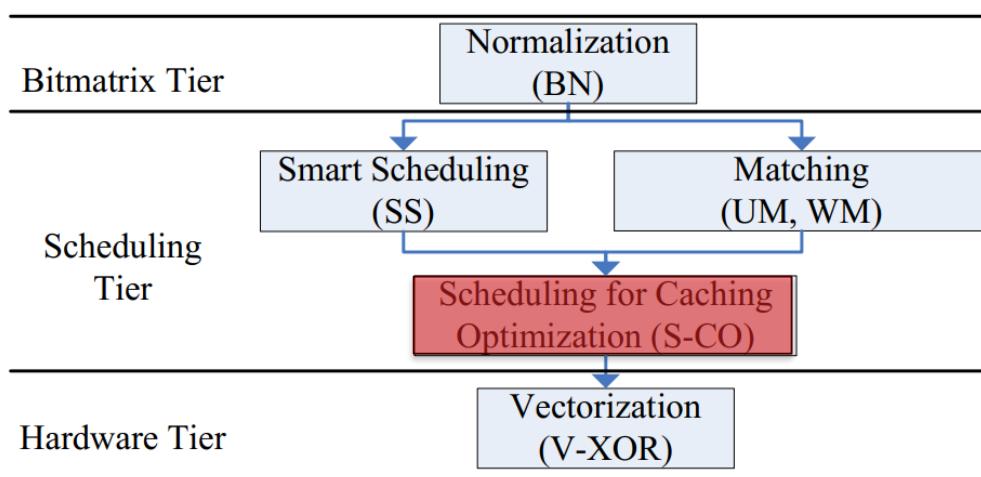
$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$



# Scheduling - Cache Optimization (S-CO)

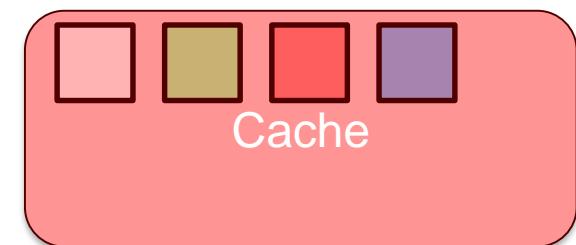
- Reduce cache miss penalty



$$p_0 = u_0 \oplus u_2 \oplus u_3 \oplus u_4$$

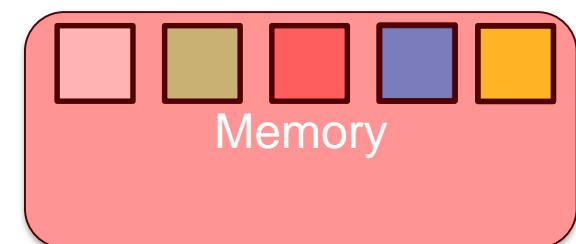
$$p_1 = u_0 \oplus u_2 \oplus u_3 \oplus u_5$$

$$u_0 \rightarrow p_0, p_1, p_2, \dots$$



$$u_1 \rightarrow p_2, p_3, p_5, \dots$$

$$u_2 \rightarrow p_0, p_1, p_4, \dots$$



$$u_3 \rightarrow p_2, p_6, p_8, \dots$$



# Question to Answer

- Most effective technique(s)?
- Utilize together?
- Components to be optimized?



# Individual Techniques

Table 1: Performance improvements by individual techniques

$(n, k, w)$	The number of XOR's in the baseline code	Reduction in the number of XOR's				Throughput increase	
		BN	SS	UM	WM	S-CO	V-XOR
(8,6,4)	104	42.31%	17.31%	31.73%	31.73%	2.14%	98.44%
(8,6,8)	362	53.31%	33.70%	34.53%	34.81%	0.52%	126.36%
(9,6,4)	152	32.89%	19.74%	32.89%	32.89%	3.54%	90.48%
(9,6,8)	549	44.63%	29.14%	38.25%	38.25%	4.51%	152.87%
(10,6,4)	200	27.50%	22.00%	36.00%	36.00%	0.60%	91.30%
(10,6,8)	736	40.90%	28.80%	38.59%	38.59%	1.23%	130.71%
(12,8,4)	256	23.44%	23.44%	35.94%	35.94%	9.21%	118.26%
(12,8,8)	1028	36.38%	24.81%	39.79%	39.79%	0.10%	156.88%
(16,10,4)	496	18.95%	21.77%	37.70%	37.70%	8.28%	138.93%
(16,10,8)	1920	30.16%	21.98%	40.89%	40.89%	5.00%	179.00%
Average over all tested cases		35.05%	24.27%	36.63%	36.66%	4.81%	130.4%



# Individual Techniques

XOR-based  
Vectorization

Table 1: Performance improvements by individual techniques

$(n, k, w)$	The number of XOR's in the baseline code	Reduction in the number of XOR's				Throughput increase	
		BN	SS	UM	WM	S-CO	V-XOR
(8,6,4)	104	42.31%	17.31%	31.73%	31.73%	2.14%	98.44%
(8,6,8)	362	53.31%	33.70%	34.53%	34.81%	0.52%	126.36%
(9,6,4)	152	32.89%	19.74%	32.89%	32.89%	3.54%	90.48%
(9,6,8)	549	44.63%	29.14%	38.25%	38.25%	4.51%	152.87%
(10,6,4)	200	27.50%	22.00%	36.00%	36.00%	0.60%	91.30%
(10,6,8)	736	40.90%	28.80%	38.59%	38.59%	1.23%	130.71%
(12,8,4)	256	23.44%	23.44%	35.94%	35.94%	9.21%	118.26%
(12,8,8)	1028	36.38%	24.81%	39.79%	39.79%	0.10%	156.88%
(16,10,4)	496	18.95%	21.77%	37.70%	37.70%	8.28%	138.93%
(16,10,8)	1920	30.16%	21.98%	40.89%	40.89%	5.00%	179.00%
Average over all tested cases		35.05%	24.27%	36.63%	36.66%	4.81%	130.4%



# Individual Techniques

XOR-based  
Vectorization

Table 1: Performance improvements by individual techniques

$(n, k, w)$	The number of XOR's in the baseline code	Reduction in the number of XOR's				Throughput increase	
		BN	SS	UM	WM	S-CO	V-XOR
(8,6,4)	104	42.31%	17.31%	31.73%	31.73%	2.14%	98.44%
(8,6,8)	362	53.31%	33.70%	34.53%	34.81%	0.52%	126.36%
(9,6,4)	152	32.89%	19.74%	32.89%	32.89%	3.54%	90.48%
(9,6,8)	549	44.63%	29.14%	38.25%	38.25%	4.51%	152.87%
(10,6,4)	200	27.50%	22.00%	36.00%	36.00%	0.60%	91.30%
(10,6,8)	736	40.90%	28.80%	38.59%	38.59%	1.23%	130.71%
(12,8,4)	256	23.44%	23.44%	35.94%	35.94%	9.21%	118.26%
(12,8,8)	1028	36.38%	24.81%	39.79%	39.79%	0.10%	156.88%
(16,10,4)	496	18.95%	21.77%	37.70%	37.70%	8.28%	138.93%
(16,10,8)	1920	30.16%	21.98%	40.89%	40.89%	5.00%	179.00%
Average over all tested cases		35.05%	24.27%	36.63%	36.66%	4.81%	130.4%

Bitmatrix  
Normalization

Unweighted  
Matching

Weighted  
Matching



# Individual Techniques

XOR-based  
Vectorization

Table 1: Performance improvements by individual techniques

$(n, k, w)$	The number of XOR's in the baseline code	Reduction in the number of XOR's				Throughput increase	
		BN	SS	UM	WM	S-CO	V-XOR
(8,6,4)	104	42.31%	17.31%	31.73%	31.73%	2.14%	98.44%
(8,6,8)	362	53.31%	33.70%	34.53%	34.81%	0.52%	126.36%
(9,6,4)	152	32.89%	19.74%	32.89%	32.89%	3.54%	90.48%
(9,6,8)	549	44.63%	29.14%	38.25%	38.25%	4.51%	152.87%
(10,6,4)	200	27.50%	22.00%	36.00%	36.00%	0.60%	91.30%
(10,6,8)	736	40.90%	28.80%	38.59%	38.59%	1.23%	130.71%
(12,8,4)	256	23.44%	23.44%	35.94%	35.94%	9.21%	118.26%
(12,8,8)	1028	36.38%	24.81%	39.79%	39.79%	0.10%	156.88%
(16,10,4)	496	18.95%	21.77%	37.70%	37.70%	8.28%	138.93%
(16,10,8)	1920	30.16%	21.98%	40.89%	40.89%	5.00%	179.00%
Average over all tested cases		35.05%	24.27%	36.63%	36.66%	4.81%	130.4%

Bitmatrix  
Normalization

Smart  
Scheduling

Unweighted  
Matching

Weighted  
Matching



# Individual Techniques

Table 1: Performance improvements by individual techniques

$(n, k, w)$	The number of XOR's in the baseline code	Reduction in the number of XOR's				Throughput increase	
		BN	SS	UM	WM	S-CO	V-XOR
(8,6,4)	104	42.31%	17.31%	31.73%	31.73%	2.14%	98.44%
(8,6,8)	362	53.31%	33.70%	34.53%	34.81%	0.52%	126.36%
(9,6,4)	152	32.89%	19.74%	32.89%	32.89%	3.54%	90.48%
(9,6,8)	549	44.63%	29.14%	38.25%	38.25%	4.51%	152.87%
(10,6,4)	200	27.50%	22.00%	36.00%	36.00%	0.60%	91.30%
(10,6,8)	736	40.90%	28.80%	38.59%	38.59%	1.23%	130.71%
(12,8,4)	256	23.44%	23.44%	35.94%	35.94%	9.21%	118.26%
(12,8,8)	1028	36.38%	24.81%	39.79%	39.79%	0.10%	156.88%
(16,10,4)	496	18.95%	21.77%	37.70%	37.70%	8.28%	138.93%
(16,10,8)	1920	30.16%	21.98%	40.89%	40.89%	5.00%	179.00%

Average over all tested cases		35.05%	24.27%	36.63%	36.66%	4.81%	130.4%
----------------------------------	--	--------	--------	--------	--------	-------	--------

Bitmatrix  
Normalization

Smart  
Scheduling

Unweighted  
Matching

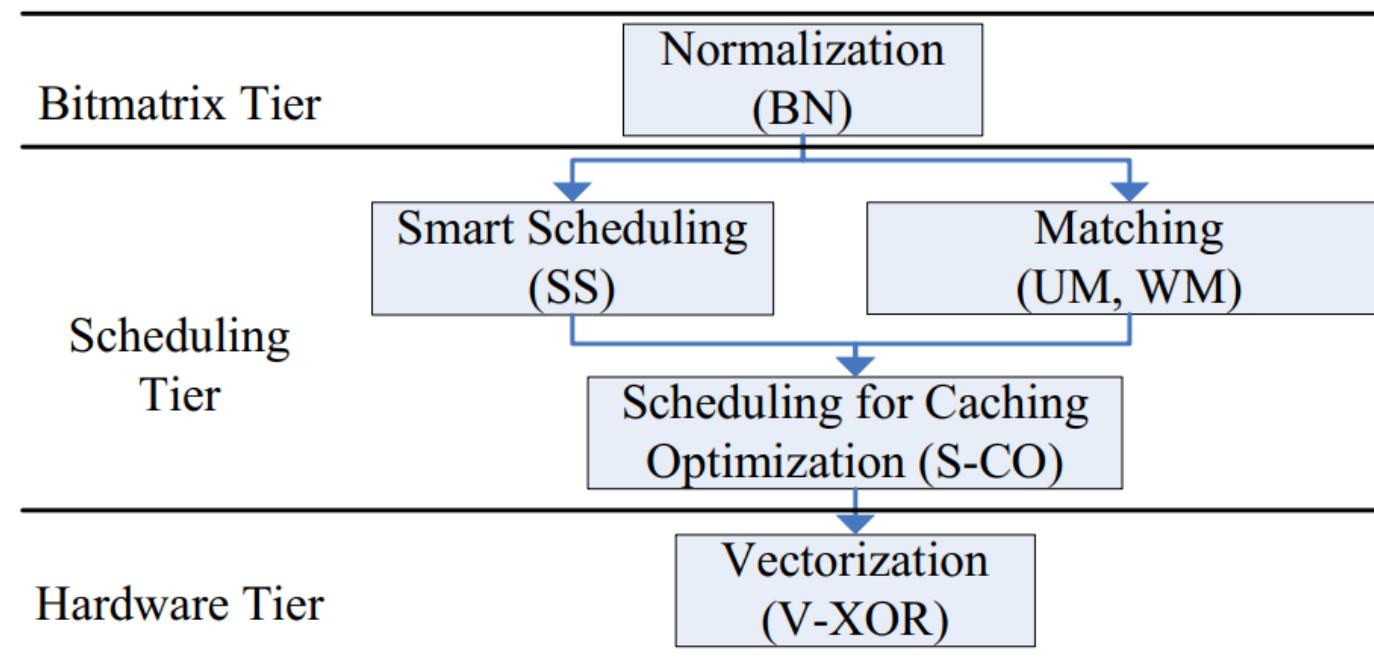
Weighted  
Matching

Cache  
Optimization

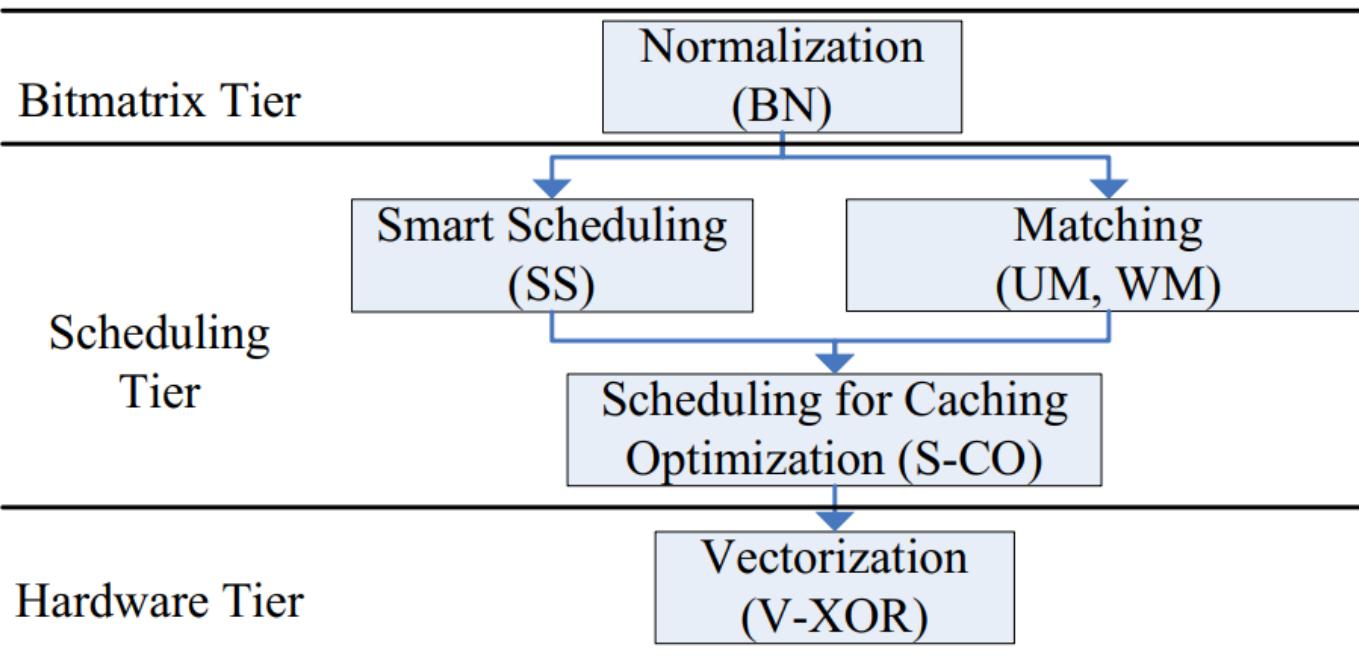
XOR-based  
Vectorization



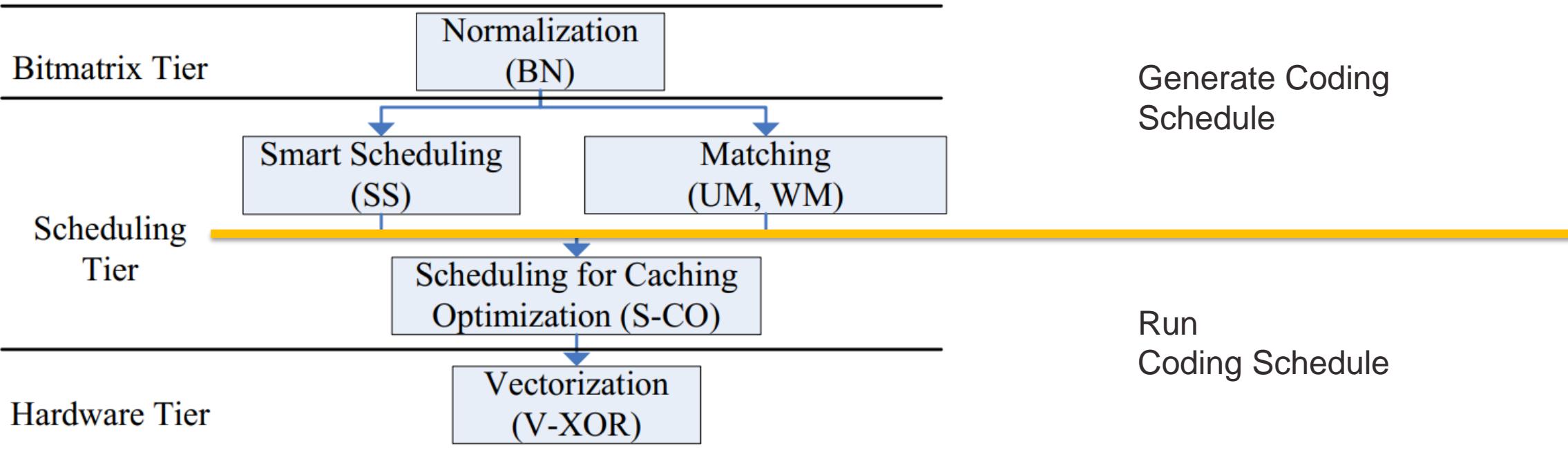
# Optimization Tiers



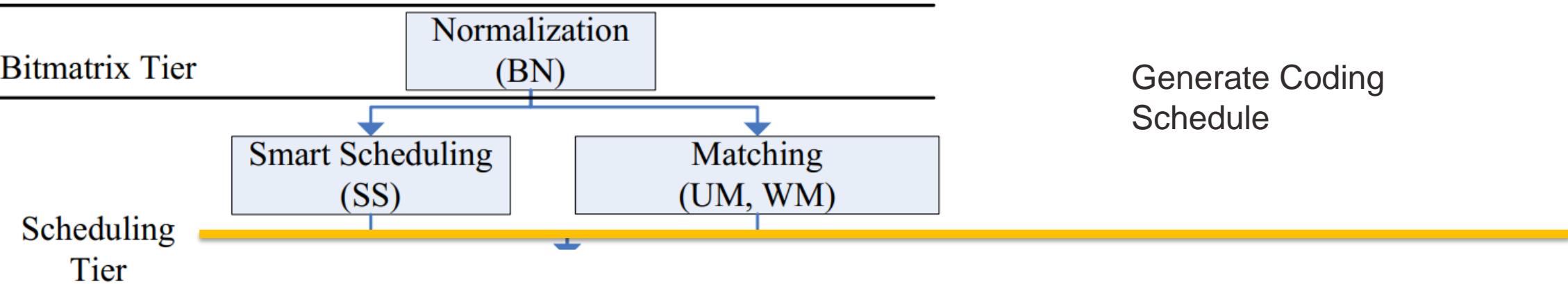
# Optimization Tiers



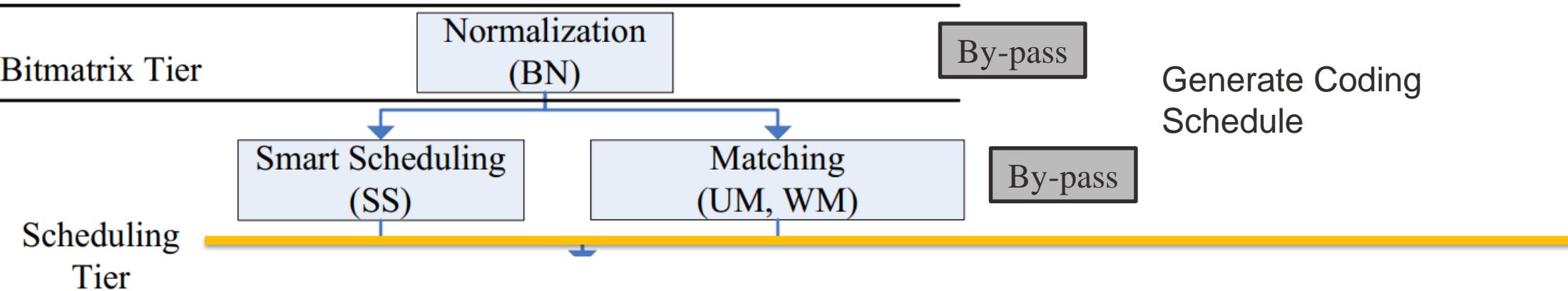
# Optimization Tiers



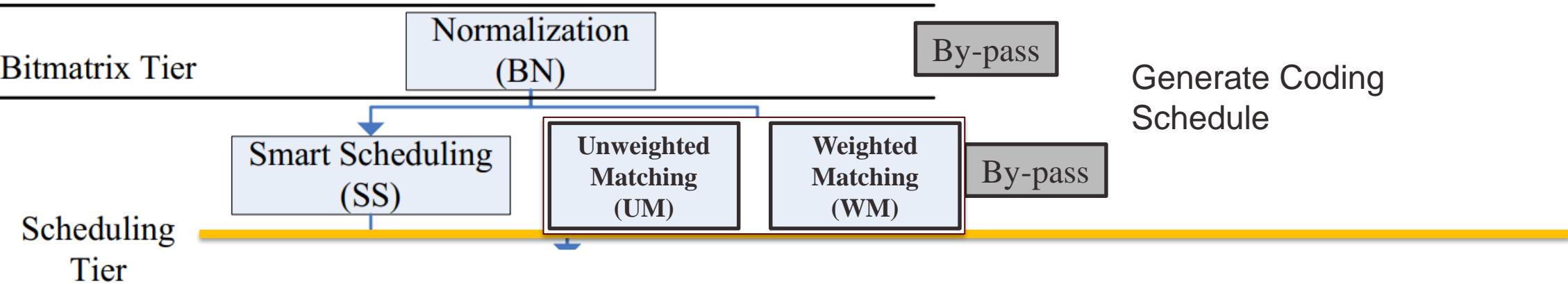
# Optimization Tiers



# Optimization Tiers



# Optimization Tiers

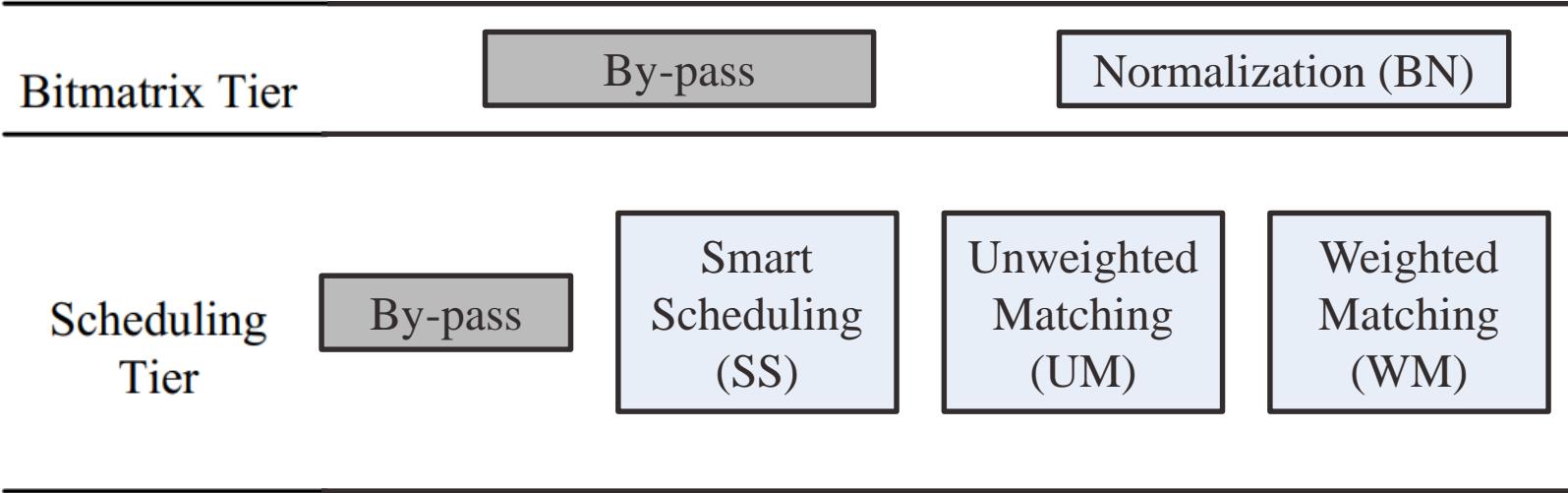


# Question to Answer

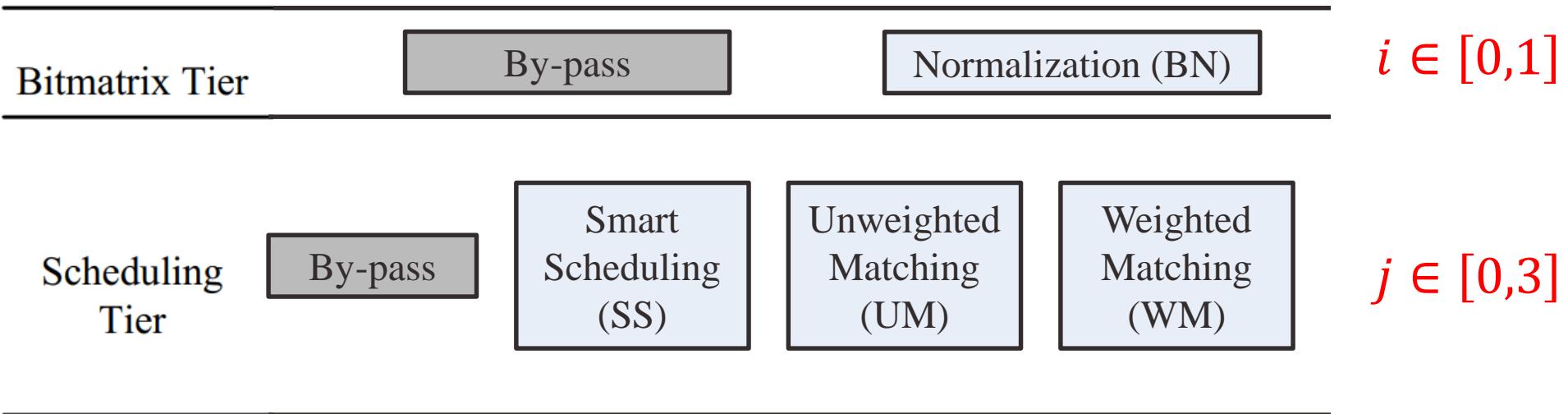
- Most effective technique(s)?
- Utilize together?
- Components to be optimized?



# Combinations (i,j)-strategy

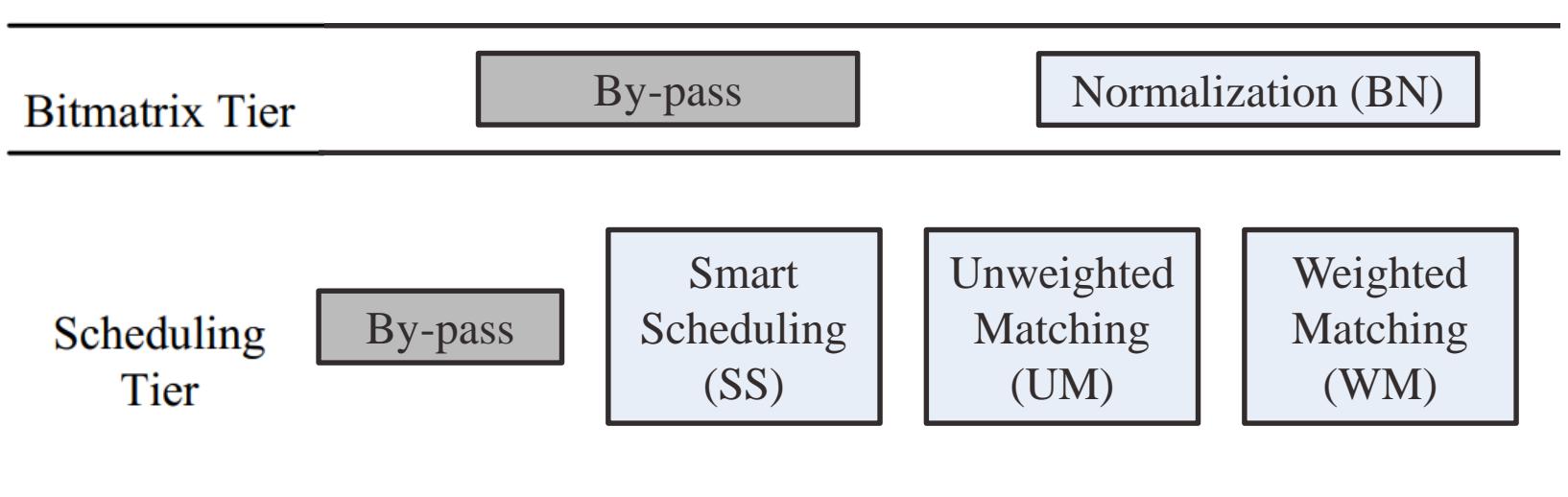


# Combinations (i,j)-strategy

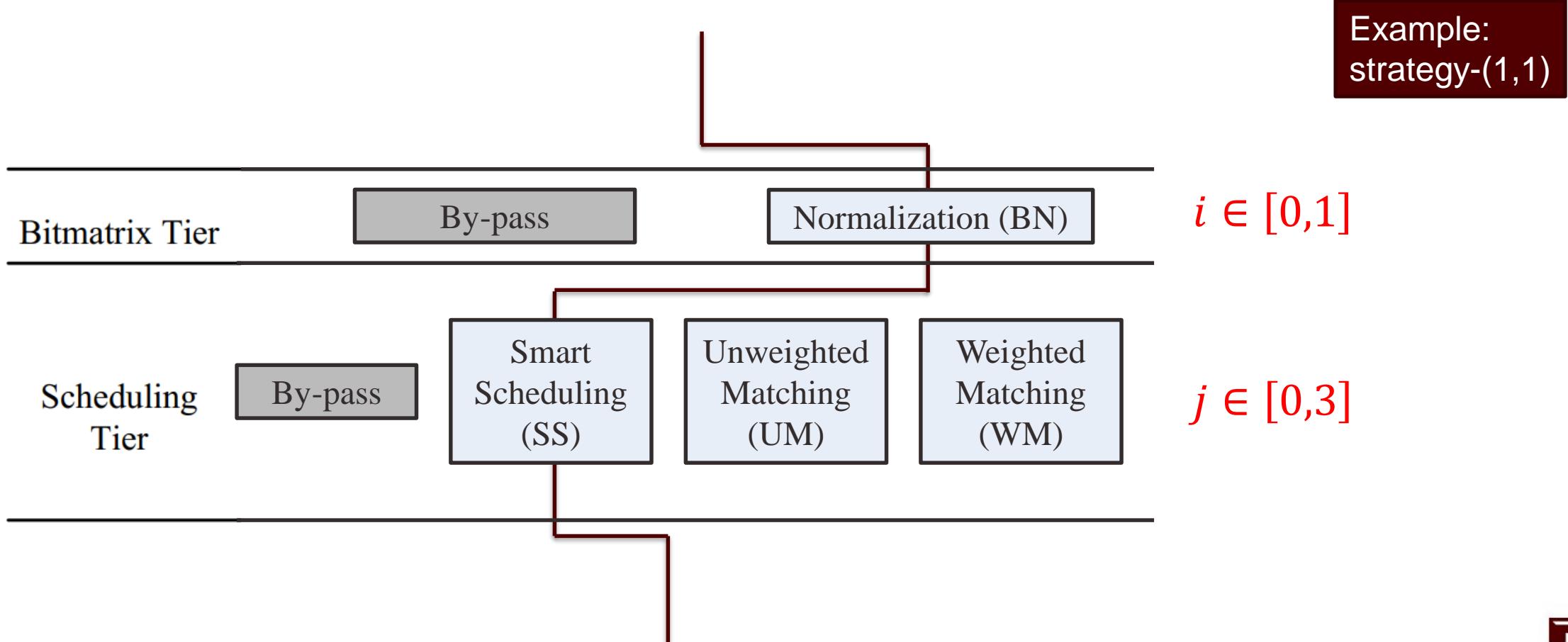


# Combinations (i,j)-strategy

Example:  
strategy-(1,1)



# Combinations (i,j)-strategy



# Contents

- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- Proposed Coding Procedure and Evaluation
- Conclusion



# Question to Answer

- Most effective technique(s)?
- Utilize together?
- Components to be optimized?



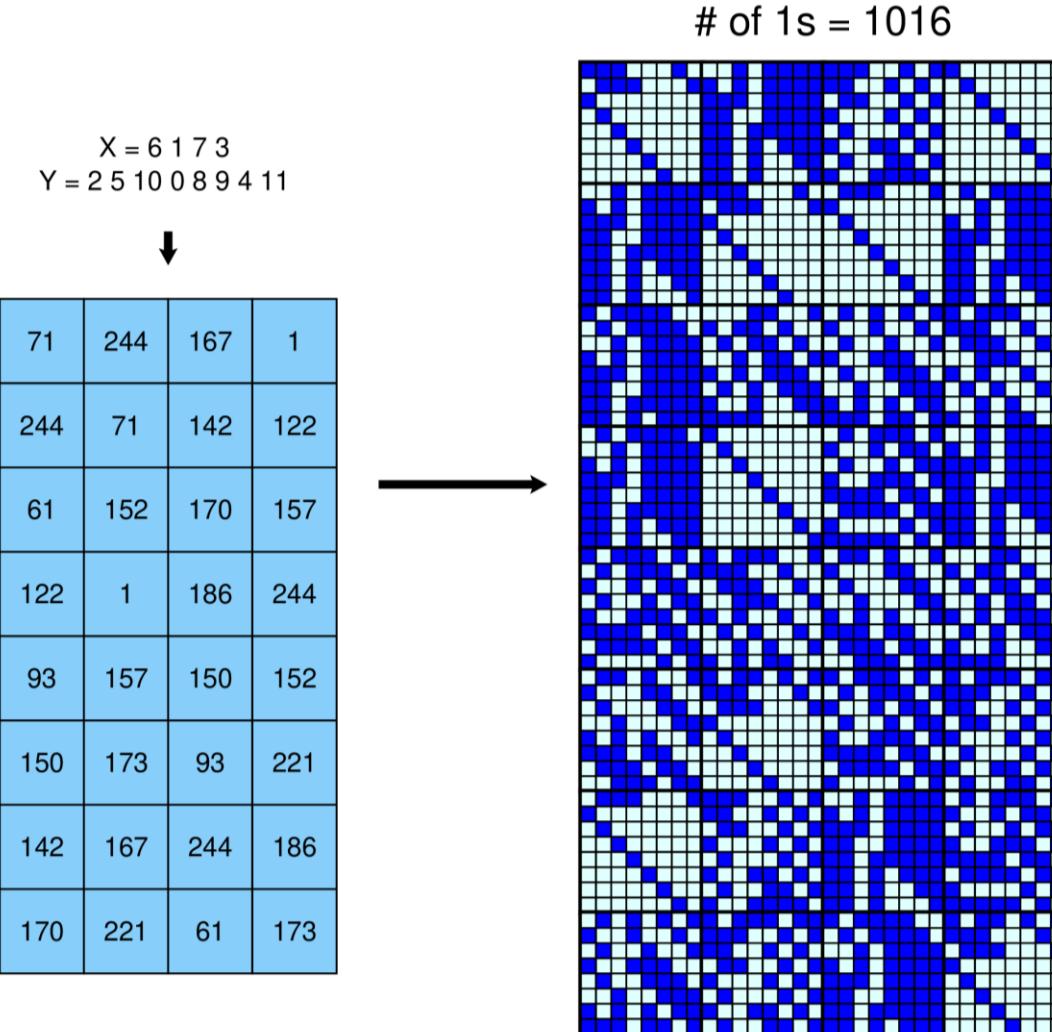
# Choice of Cauchy Matrix

- Different  $X, Y$  array yields different Cauchy Matrix



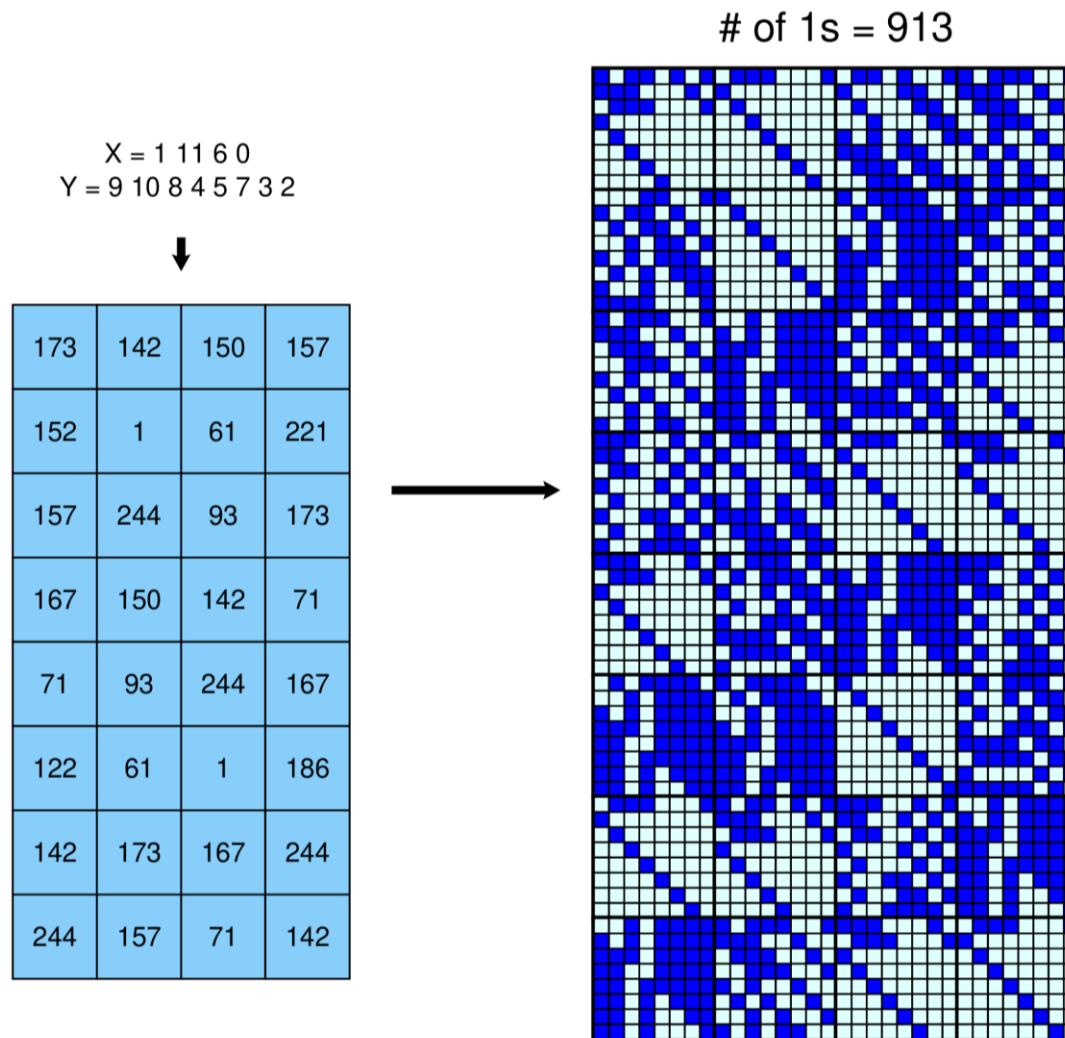
# Choice of Cauchy Matrix

- Different  $X, Y$  array yields different Cauchy Matrix



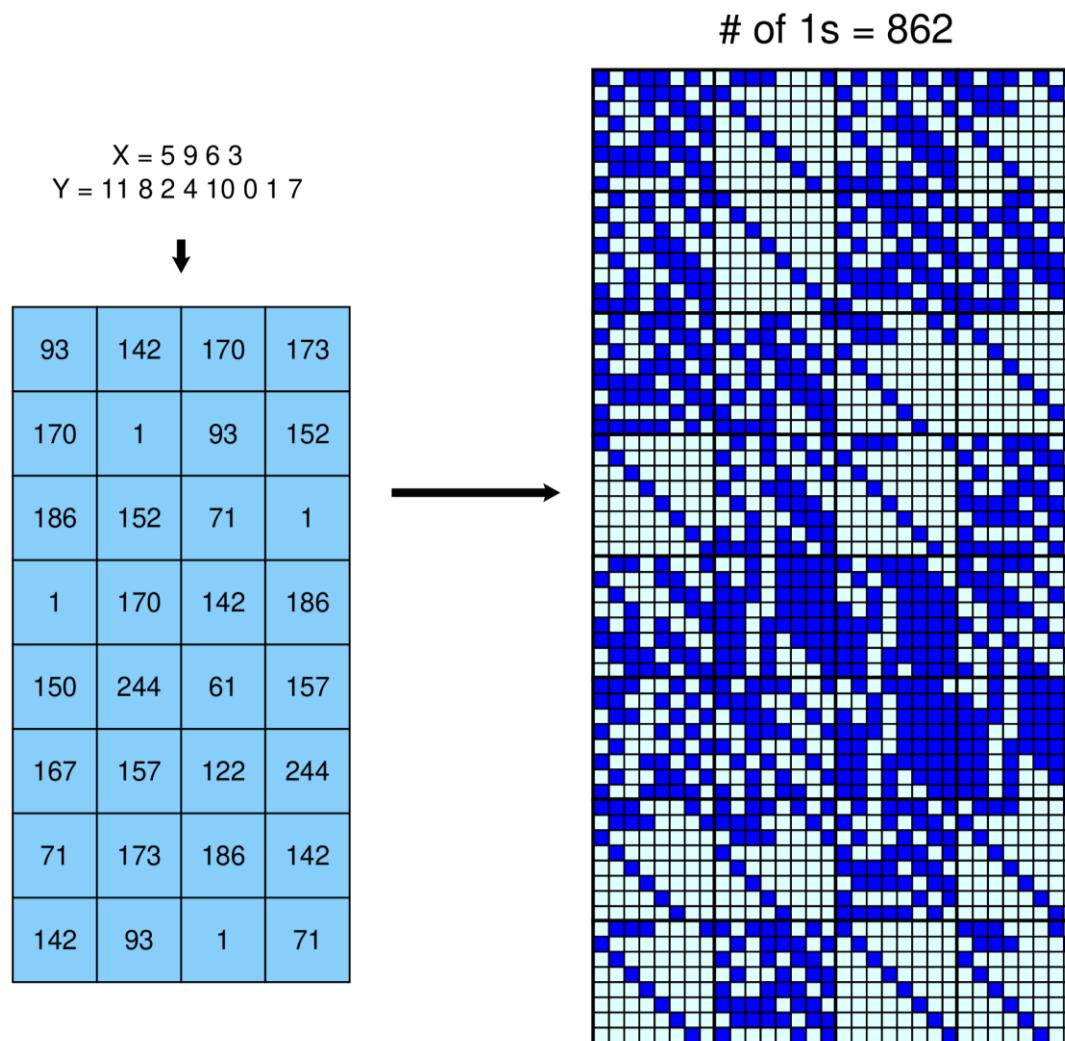
# Choice of Cauchy Matrix

- Different  $X, Y$  array yields different Cauchy Matrix



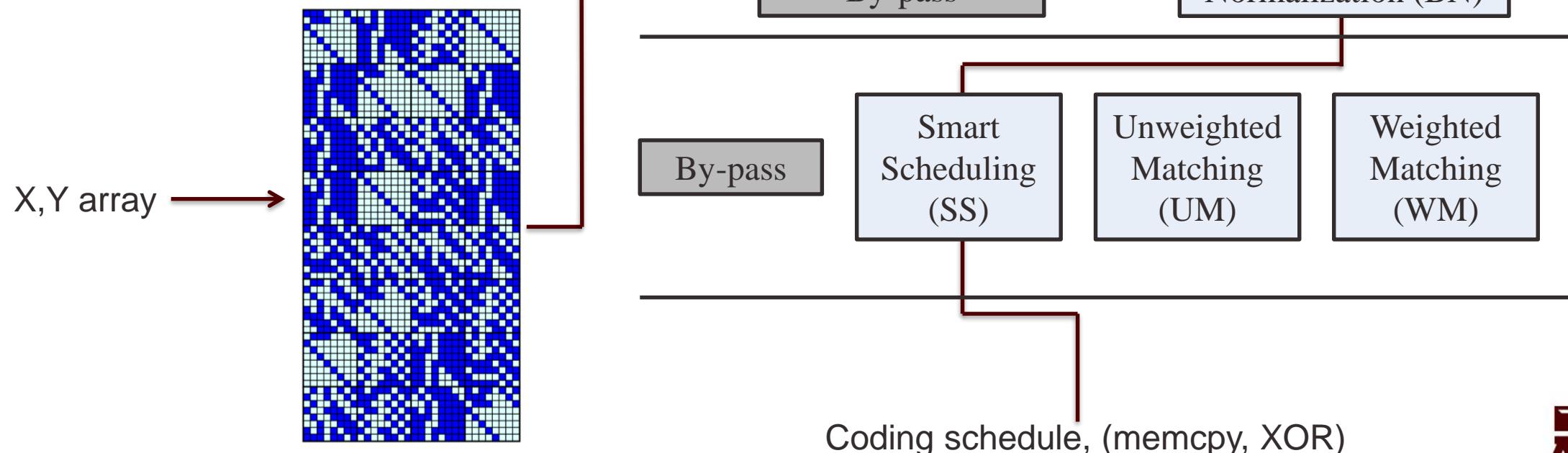
# Choice of Cauchy Matrix

- Different  $X, Y$  array yields different Cauchy Matrix



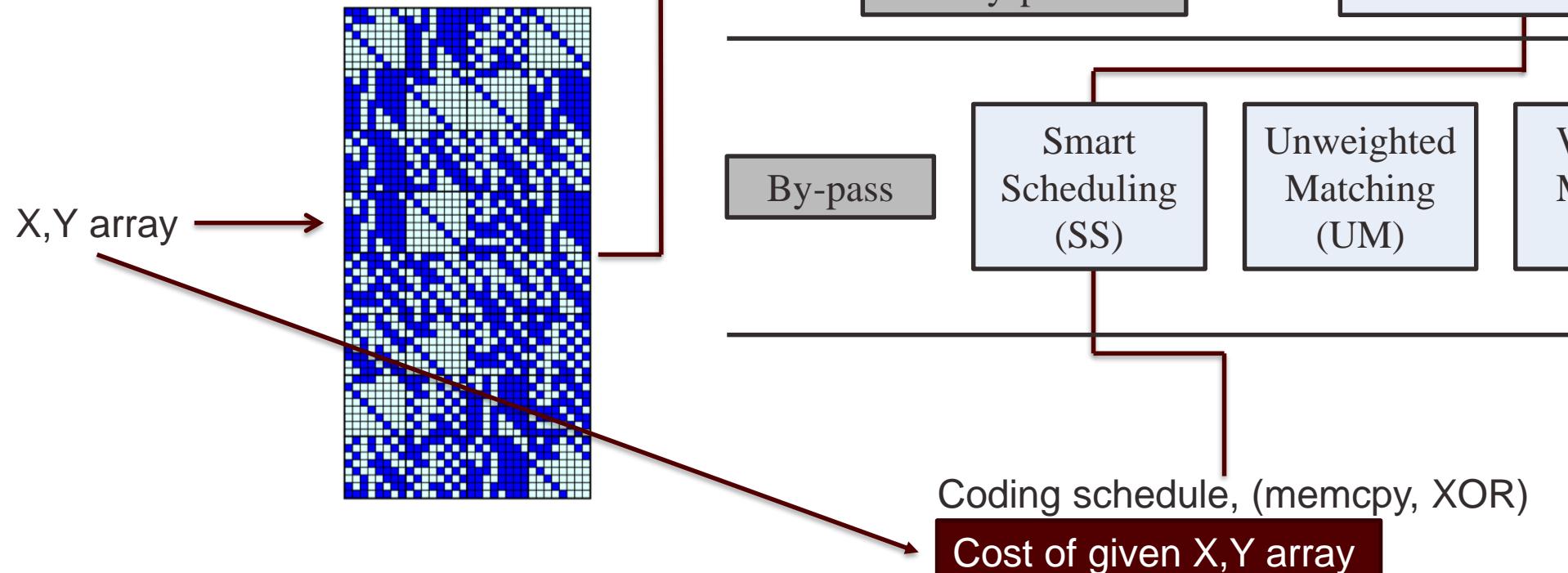
# Choice of Cauchy Matrix

- Different Cauchy Matrix affect coding chain



# Choice of Cauchy Matrix

- Different Cauchy Matrix affect coding chain



# Bitmatrix Optimization

- Individual bitmatrix optimization for all  $(i,j)$ -strategies
- Simulated Annealing and Genetic Algorithm
- # of operations as cost



# [no\_opt by SA by GA ], # of ops in schedule

$(i, j)$	$(n, k, w) = (8, 6, 4)$			$(n, k, w) = (9, 6, 4)$			$(n, k, w) = (10, 6, 4)$			$(n, k, w) = (12, 8, 4)$			$(n, k, w) = (16, 10, 4)$		
(0,0)	112	77	77	164	122	117	216	173	162	272	232	232	520	462	464
(0,1)	94	70	68	134	102	100	172	137	134	212	190	188	412	368	368
(0,2)	90	72	68	127	103	101	164	134	132	204	186	180	376	344	345
(0,3)	90	72	68	127	102	101	164	132	132	204	184	182	376	343	345
(1,0)	68	58	58	114	99	98	161	142	141	212	198	198	426	419	419
(1,1)	64	58	58	99	90	89	138	120	120	189	174	171	365	352	352
(1,2)	64	58	57	98	87	87	133	118	118	176	165	164	326	317	316
(1,3)	64	58	57	98	90	87	132	118	118	175	164	164	326	316	316

$(i, j)$	$(n, k, w) = (8, 6, 8)$			$(n, k, w) = (9, 6, 8)$			$(n, k, w) = (10, 6, 8)$			$(n, k, w) = (12, 8, 8)$			$(n, k, w) = (16, 10, 8)$		
(0,0)	378	291	247	573	467	425	768	611	582	1060	880	841	1968	1685	1681
(0,1)	256	234	225	413	393	349	556	518	479	805	740	684	1546	1432	1374
(0,2)	286	227	217	408	346	321	532	474	450	726	636	591	1304	1180	1170
(0,3)	286	216	197	408	357	329	532	460	450	726	636	627	1304	1180	1170
(1,0)	185	151	142	328	286	262	467	434	419	686	613	563	1389	1293	1246
(1,1)	164	155	133	285	270	230	411	383	371	593	557	544	1264	1184	1133
(1,2)	167	143	134	272	244	225	377	343	335	520	487	462	998	951	922
(1,3)	167	144	130	273	249	233	377	337	337	520	473	467	995	941	932

[no\_opt by SA by GA ], # of ops in schedule

$(i, j)$	$(n, k, w) = (8, 6, 4)$			$(n, k, w) = (9, 6, 4)$			$(n, k, w) = (10, 6, 4)$			$(n, k, w) = (12, 8, 4)$			$(n, k, w) = (16, 10, 4)$		
(0,0)	112	77	77	164	122	117	216	173	162	272	232	232	520	462	464
(0,1)	94	70	68	134	102	100	172	137	134	212	190	188	412	368	368
(0,2)	90	72	68	127	103	101	164	134	132	204	186	180	376	344	345
(0,3)	90	72	68	127	102	101	164	132	132	204	184	182	376	343	345
(1,0)	68	58	58	114	99	98	161	142	141	212	198	198	426	419	419
(1,1)	64	58	58	99	90	89	138	120	120	189	174	171	365	352	352
(1,2)	64	58	57	98	87	87	133	118	118	176	165	164	326	317	316
(1,3)	64	58	57	98	90	87	132	118	118	175	164	164	326	316	316

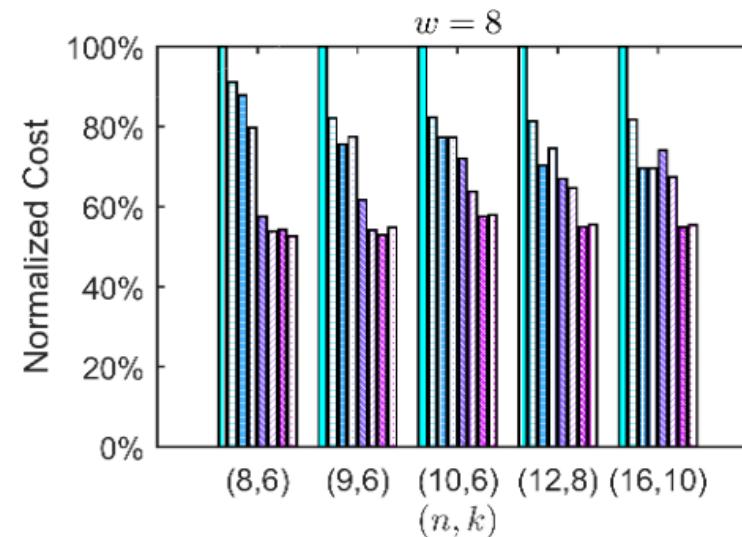
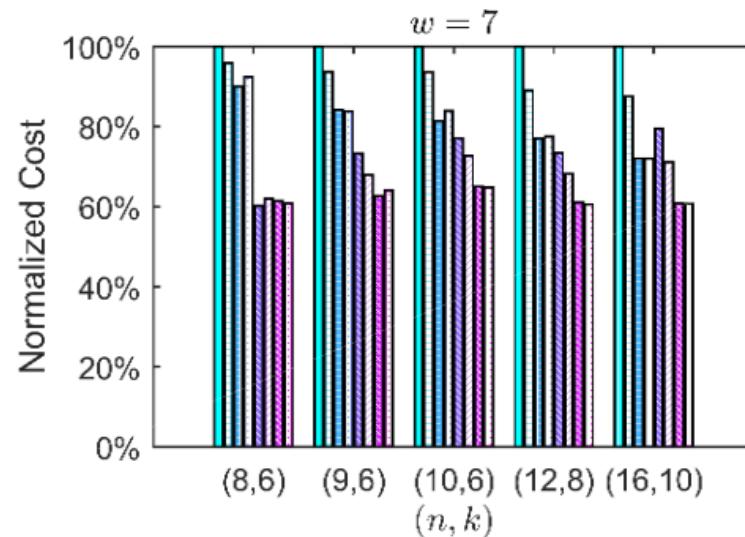
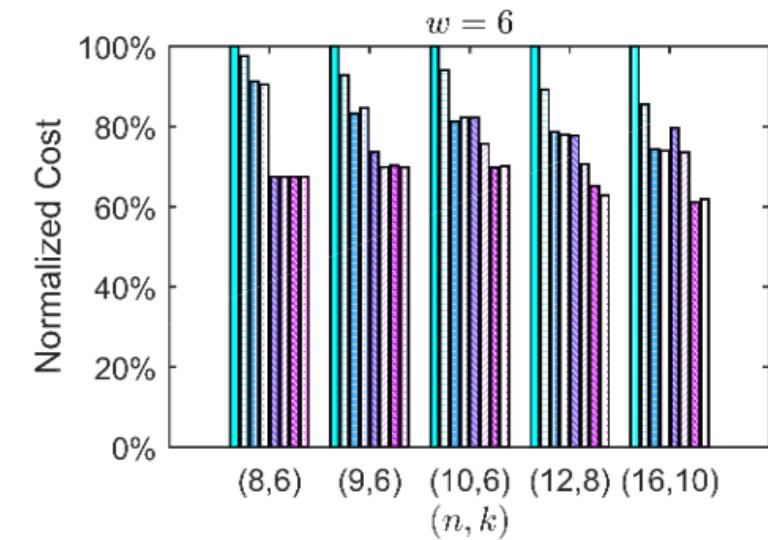
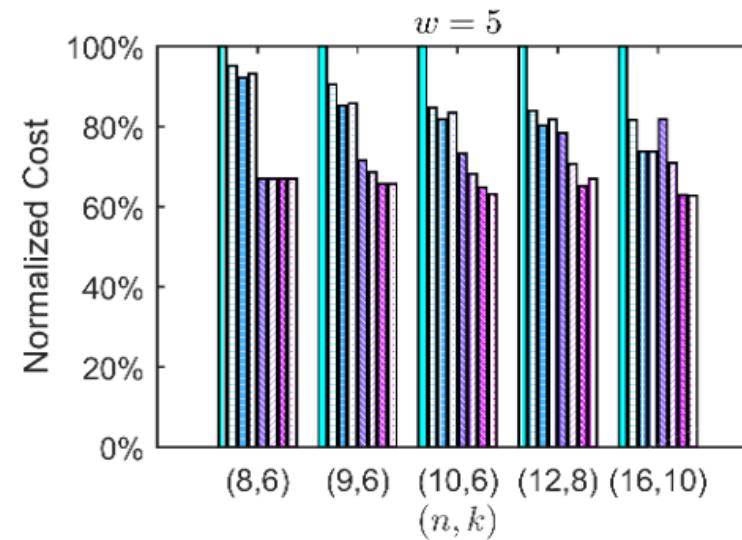
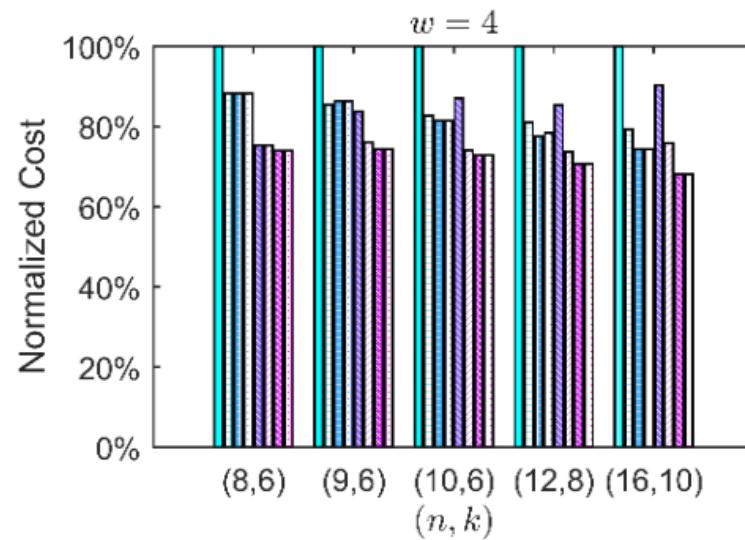
$(i, j)$	$(n, k, w) = (8, 6, 8)$			$(n, k, w) = (9, 6, 8)$			$(n, k, w) = (10, 6, 8)$			$(n, k, w) = (12, 8, 8)$			$(n, k, w) = (16, 10, 8)$		
(0,0)	378	291	247	573	467	425	768	611	582	1060	880	841	1968	1685	1681
(0,1)	256	234	225	413	393	349	556	518	479	805	740	684	1546	1432	1374
(0,2)	286	227	217	408	346	321	532	474	450	726	636	591	1304	1180	1170
(0,3)	286	216	197	408	357	329	532	460	450	726	636	627	1304	1180	1170
(1,0)	185	151	142	328	286	262	467	434	419	686	613	563	1389	1293	1246
(1,1)	164	155	133	285	270	230	411	383	371	593	557	544	1264	1184	1133
(1,2)	167	143	134	272	244	225	377	343	335	520	487	462	998	951	922
(1,3)	167	144	130	273	249	233	377	337	337	520	473	467	995	941	932

[no\_opt by SA by GA ], # of ops in schedule

$(i, j)$	$(n, k, w) = (8, 6, 4)$			$(n, k, w) = (9, 6, 4)$			$(n, k, w) = (10, 6, 4)$			$(n, k, w) = (12, 8, 4)$			$(n, k, w) = (16, 10, 4)$		
(0,0)	112	77	77	164	122	117	216	173	162	272	232	232	520	462	464
(0,1)	94	70	68	134	102	100	172	137	134	212	190	188	412	368	368
(0,2)	90	72	68	127	103	101	164	134	132	204	186	180	376	344	345
(0,3)	90	72	68	127	102	101	164	132	132	204	184	182	376	343	345
(1,0)	68	58	58	114	99	98	161	142	141	212	198	198	426	419	419
(1,1)	64	58	58	99	90	89	138	120	120	189	174	171	365	352	352
(1,2)	64	58	57	98	87	87	133	118	118	176	165	164	326	317	316
(1,3)	64	58	57	98	90	87	132	118	118	175	164	164	326	316	316

$(i, j)$	$(n, k, w) = (8, 6, 8)$			$(n, k, w) = (9, 6, 8)$			$(n, k, w) = (10, 6, 8)$			$(n, k, w) = (12, 8, 8)$			$(n, k, w) = (16, 10, 8)$		
(0,0)	378	291	247	573	467	425	768	611	582	1060	880	841	1968	1685	1681
(0,1)	256	234	225	413	393	349	556	518	479	805	740	684	1546	1432	1374
(0,2)	286	227	217	408	346	321	532	474	450	726	636	591	1304	1180	1170
(0,3)	286	216	197	408	357	329	532	460	450	726	636	627	1304	1180	1170
(1,0)	185	151	142	328	286	262	467	434	419	686	613	563	1389	1293	1246
(1,1)	164	155	133	285	270	230	411	383	371	593	557	544	1264	1184	1133
(1,2)	167	143	134	272	244	225	377	343	335	520	487	462	998	951	922
(1,3)	167	144	130	273	249	233	377	337	337	520	473	467	995	941	932

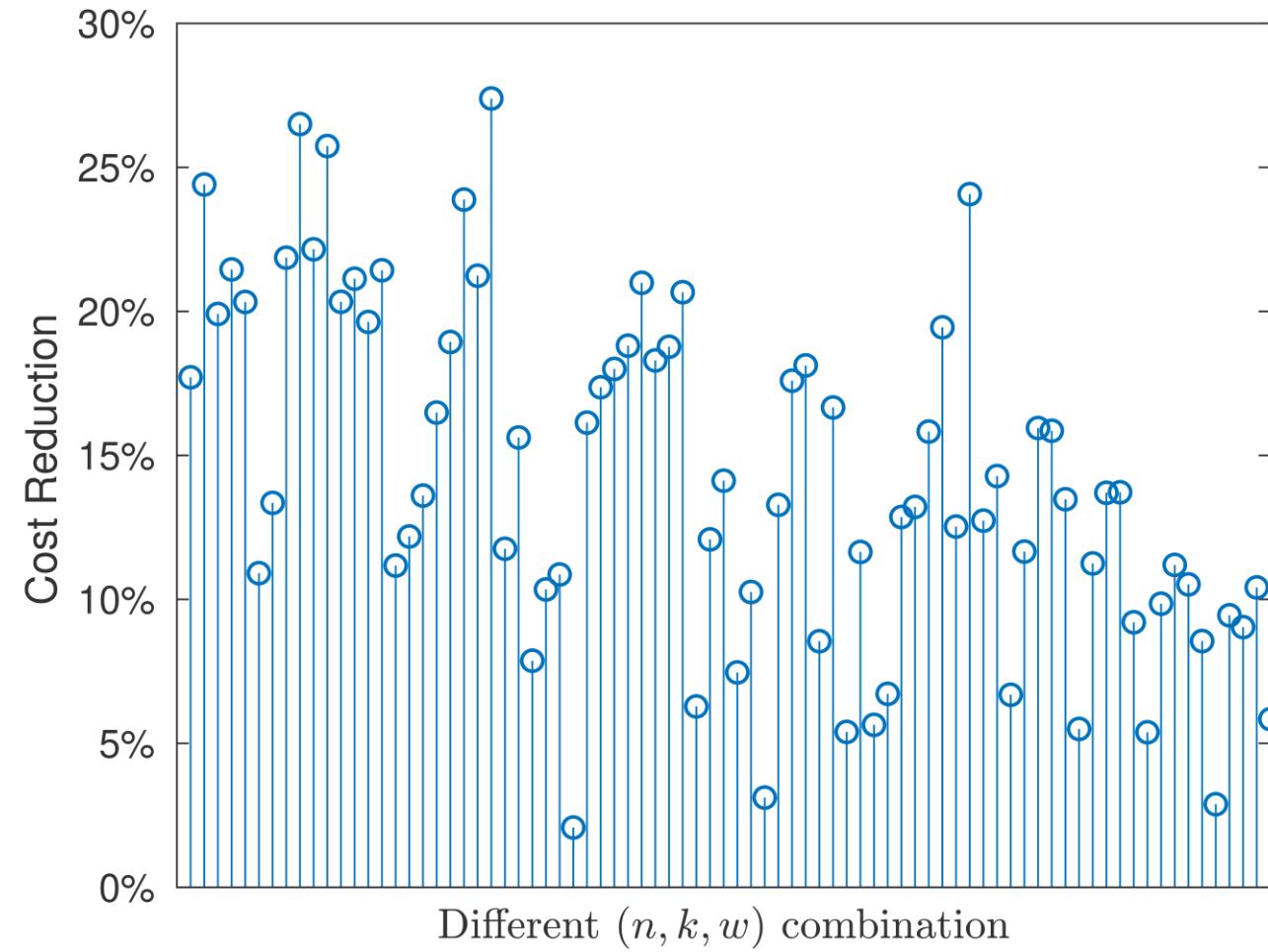
# Strategies-(i,j) with Optimized Bitmatrices



strategy-(1,3) is  
the best strategy



# Optimized Bitmatrix Reduced Costs



# Cost Function Improvement

- Memcpy is 1.5x faster than XOR
- Cost functions:
  - Total # of XORs in schedule
  - Total # of ops in schedule
  - Weighted sum of ops in schedule

Table 4: Encoding throughput (GB/s) using bitmatrices obtained by the genetic algorithm under different cost functions

$(n, k, w)$	Cost Function		
	# of XOR	# of XOR and copying	Weighted
(8,6,4)	4.64	4.66	<b>4.68</b>
(8,6,8)	4.30	<b>4.35</b>	4.32
(9,6,4)	3.72	3.73	<b>3.80</b>
(9,6,8)	3.28	3.28	<b>3.44</b>
(10,6,4)	2.33	2.51	<b>2.52</b>
(10,6,8)	1.99	1.97	<b>2.11</b>
(12,8,4)	2.96	3.11	<b>3.16</b>
(12,8,8)	2.54	2.56	<b>2.58</b>
(16,10,4)	2.29	2.29	<b>2.32</b>
(16,10,8)	1.71	1.72	<b>1.74</b>

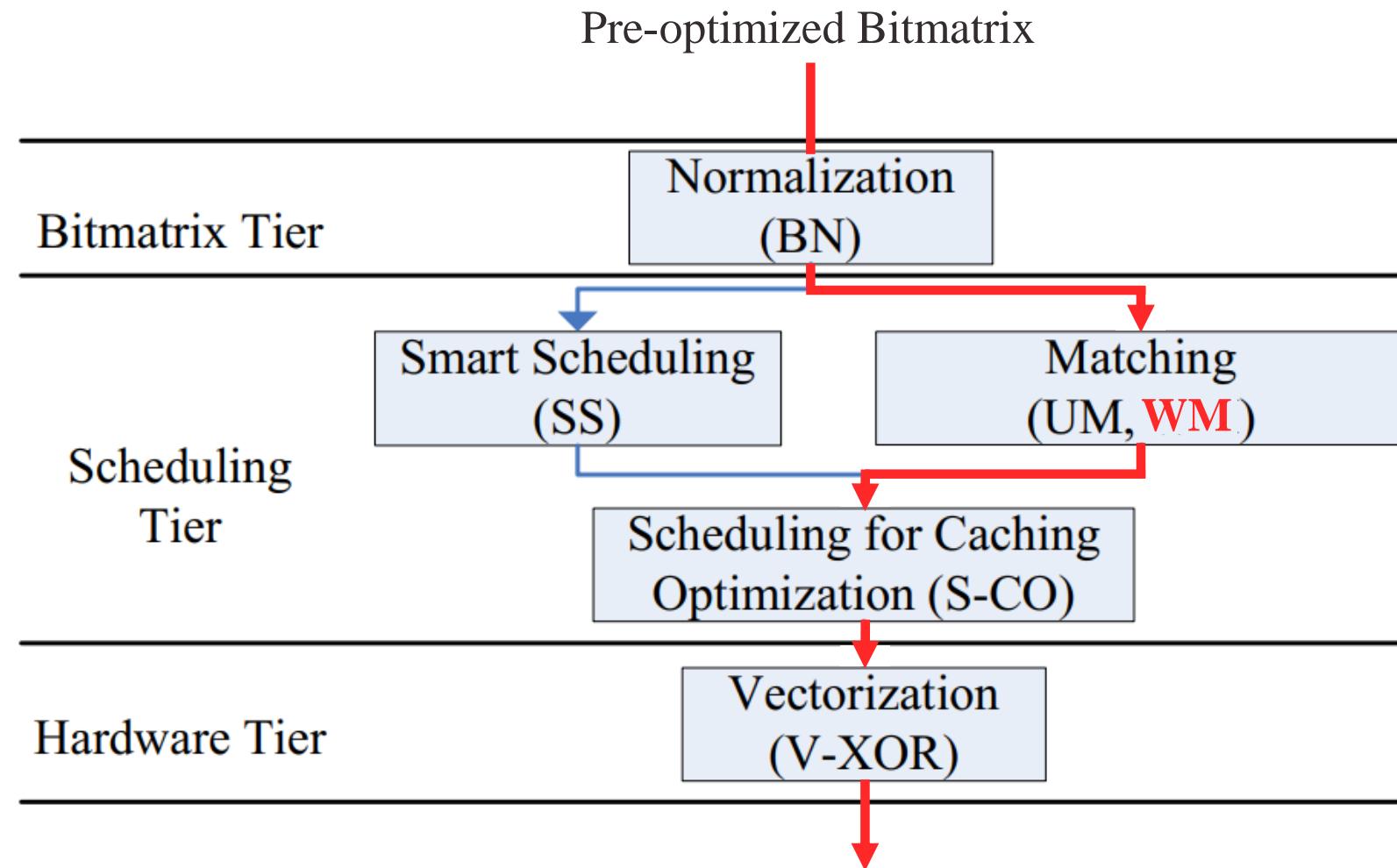


# Contents

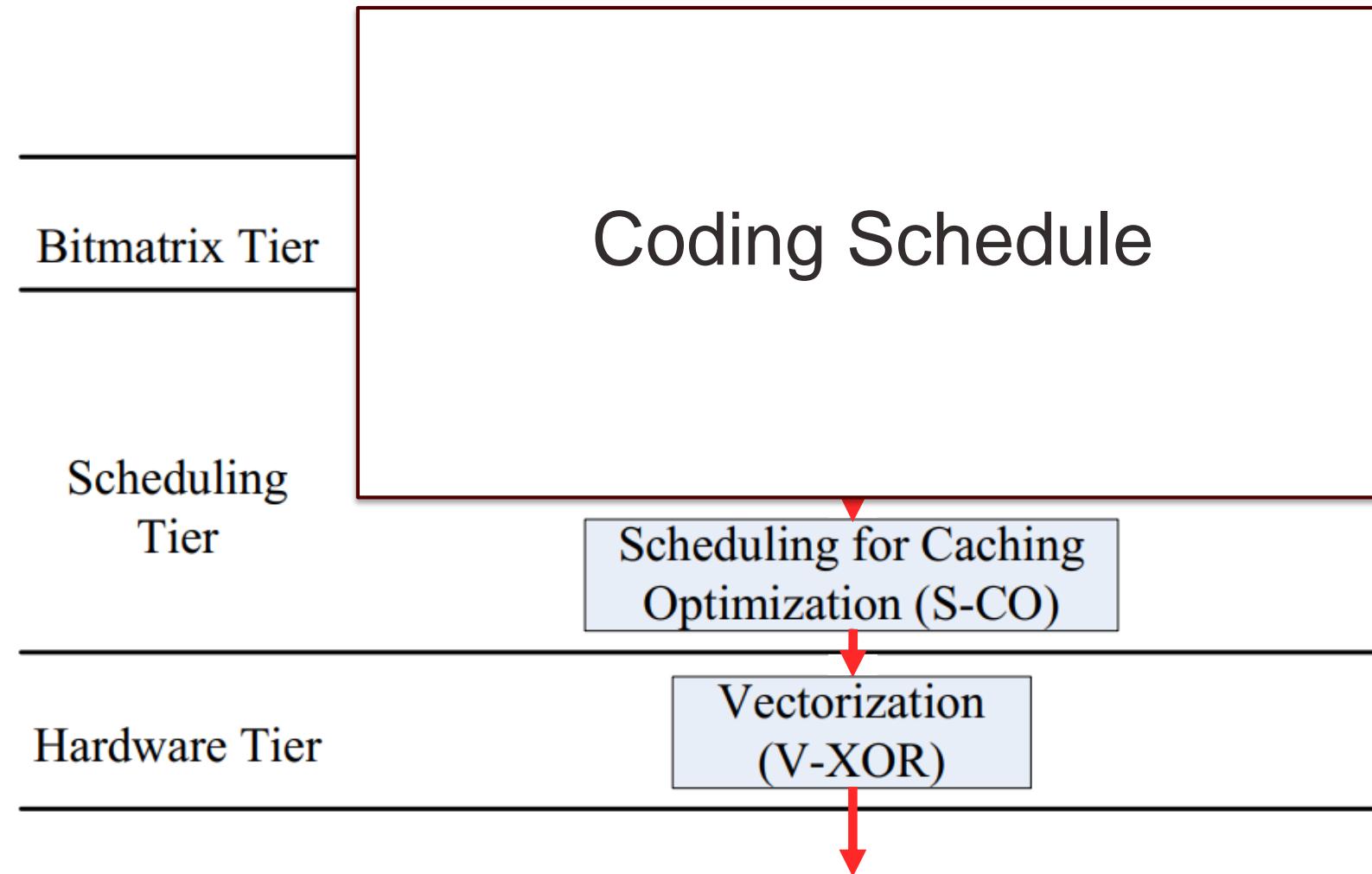
- Motivation
- Background and Review
- Evaluating Individual Techniques
- Find the Best Strategy under Optimized Bitmatrix
- **Proposed Coding Procedure and Evaluation**
- Conclusion



# Proposed Coding Procedure



# Proposed Coding Procedure



# Testing Setup

- Ryzen 1700X @ 3.4Ghz (8C/16T)
- 16GB DDR4
- Ubuntu 18.04 64-bit, GCC 7.3.0
- Jerasure 1.2A/ Jerasure 2.0:
  - XOR-based Cauchy RS code (BN, SS applied)
  - GF-based RS code
  - Raid-6
- Other codes:
  - EVENODD
  - RDP
  - STAR



# Encoding v.s. Efficient RS/CRS code

Table 5: Encoding throughput (GB/s) for methods that allow general  $(n, k)$  parameters and  $w = 8$

$(n, k)$	Proposed	Vectorized XOR-based CRS code	Vectorized GF-based RS code [16]
(7,5)	4.64	<b>4.73</b>	2.52
(8,6)	5.21	<b>5.22</b>	2.70
(9,7)	5.32	<b>5.45</b>	2.74
(10,8)	5.36	<b>5.59</b>	2.77
(12,10)	5.72	<b>5.88</b>	2.81
(8,5)	<b>3.19</b>	2.75	1.76
(9,6)	<b>3.49</b>	2.84	1.77
(10,7)	<b>3.67</b>	2.79	1.80
(11,8)	<b>3.72</b>	2.92	1.82
(13,10)	<b>3.82</b>	3.10	1.84
(10,6)	<b>2.55</b>	2.15	1.31
(11,7)	<b>2.75</b>	2.17	1.32
(12,8)	<b>2.86</b>	2.20	1.35
(14,10)	<b>2.86</b>	2.19	1.40
(15,10)	<b>2.30</b>	1.79	1.11
(16,10)	<b>1.96</b>	1.48	0.92

Outperform vectorized XOR-based CRS code for  $m > 2$



# Encoding v.s. Efficient RS/CRS code

Table 5: Encoding throughput (GB/s) for methods that allow general  $(n, k)$  parameters and  $w = 8$

$(n, k)$	Proposed	Vectorized XOR-based CRS code	Vectorized GF-based RS code [16]
(7,5)	4.64	4.73	2.52
(8,6)	5.21	5.22	2.70
(9,7)	5.32	5.45	2.74
(10,8)	5.36	5.59	2.77
(12,10)	5.72	5.88	2.81
(8,5)	3.19	2.75	1.76
(9,6)	3.49	2.84	1.77
(10,7)	3.67	2.79	1.80
(11,8)	3.72	2.92	1.82
(13,10)	3.82	3.10	1.84
(10,6)	2.55	2.15	1.31
(11,7)	2.75	2.17	1.32
(12,8)	2.86	2.20	1.35
(14,10)	2.86	2.19	1.40
(15,10)	2.30	1.79	1.11
(16,10)	1.96	1.48	0.92

Outperform vectorized XOR-based CRS code for  $m > 2$

Outperform vectorized GF-based RS code with big margin (~ 2x faster)

# Encoding v.s. Three Parities Codes

Table 6: Encoding throughputs (GB/s): Three parities

$(n, k, w)$	Proposed	Vectorized XOR-based CRS code	STAR code [8]	Quantcast QFS [13]
(8,5,4)	<b>3.59</b>	3.25	2.97	2.92
(8,5,8)	<b>3.19</b>	2.75	2.97	2.92
(9,6,4)	3.52	<b>3.72</b>	3.42	3.04
(9,6,8)	<b>3.49</b>	2.84	3.42	3.04
(10,7,4)	<b>4.15</b>	3.86	3.76	3.25
(10,7,8)	3.67	2.79	<b>3.76</b>	3.25
(11,8,4)	<b>4.36</b>	4.13	3.94	3.27
(11,8,8)	3.72	2.92	<b>3.94</b>	3.27
(13,10,4)	<b>4.51</b>	4.08	4.37	3.41
(13,10,8)	3.82	3.10	<b>4.37</b>	3.41

Similar or better performance than specially designed codes

# Encoding v.s. Two Parities Codes

Table 7: Encoding throughputs (GB/s): Two parities

$(n, k, w)$	Proposed	Vectorized XOR-based CRS code	Vectorized Raid-6	EVEN ODD [2]	RDP [6]
(7,5,4)	<b>5.16</b>	4.85	n/a	4.28	4.37
(7,5,8)	4.64	<b>4.73</b>	2.18	4.28	4.37
(8,6,4)	4.67	<b>5.22</b>	n/a	4.83	4.95
(8,6,8)	5.21	<b>5.22</b>	2.15	4.83	4.95
(9,7,4)	<b>5.77</b>	5.59	n/a	5.20	5.32
(9,7,8)	5.32	<b>5.45</b>	2.16	5.20	5.32
(10,8,4)	<b>5.90</b>	5.23	n/a	5.50	5.69
(10,8,8)	5.36	<b>5.59</b>	2.17	5.50	<b>5.69</b>
(12,10,4)	6.23	6.00	n/a	6.02	<b>6.24</b>
(12,10,8)	5.72	5.88	2.17	6.02	<b>6.24</b>

Similar or better performance than specially designed codes



# Overall Encoding Improvement

Reference codes or methods	Improvement by proposed code
General $(n, k)$ Codes	
GF-based RS code w/o vectorization	552.27%
XOR-based CRS code w/o vectorization	53.65%
Vectorized GF-based RS code [16]	99.82%
Vectorized XOR-based CRS code	14.98%
Three Parities Codes	
STAR [8]	5.59%
Quancast-QFS [13]	21.68%
Two Parities Codes	
Raid-6 w/o vectorization	206.88%
Vectorized Raid-6	142.07%
RDP [6]	5.85%
EVENODD [2]	8.79%



# Decoding



# Decoding

- Direct read in most cases



# Decoding

- Direct read in most cases
- Degraded read (decoding) only if **systematic** nodes fail



# Decoding

- Direct read in most cases
- Degraded read (decoding) only if **systematic** nodes fail
  - Single disk failure rate : 1%
  - In (10,6) erasure code storage system
    - 1 systematic node failure: ~ 0.055
    - 2 systematic node failure: ~ 0.0014
    - 3 systematic node failure: ~  $1.8 \times 10^{-5}$
    - 4 systematic node failure: ~  $1.4 \times 10^{-7}$



# Decoding

- Direct read in most cases
- Degraded read (decoding) only if **systematic** nodes fail
  - Single disk failure rate : 1%
  - In (10,6) erasure code storage system
    - 1 systematic node failure: ~ 0.055
    - 2 systematic node failure: ~ 0.0014
    - 3 systematic node failure: ~  $1.8 \times 10^{-5}$
    - 4 systematic node failure: ~  $1.4 \times 10^{-7}$
- Decoding performance should be viewed as secondary importance



# Decoding Throughput

Table 9: Decoding throughput (GB/s) for methods that allow general  $(n, k)$  parameters and  $w = 8$

$(n, k)$	Proposed	Vectorized XOR-based CRS code	Vectorized GF-based RS code [16]
(7,5)	3.87	<b>4.56</b>	2.58
(8,6)	<b>5.45</b>	4.86	2.67
(9,7)	4.46	<b>5.06</b>	2.70
(10,8)	4.89	<b>5.11</b>	2.75
(12,10)	4.45	<b>5.52</b>	2.79
(8,5)	<b>3.04</b>	2.11	1.71
(9,6)	<b>2.94</b>	2.20	1.74
(10,7)	<b>3.28</b>	2.29	1.76
(11,8)	<b>3.08</b>	2.31	1.71
(13,10)	<b>3.21</b>	2.37	1.88
(10,6)	<b>2.38</b>	1.80	1.31
(11,7)	<b>2.35</b>	1.85	1.32
(12,8)	<b>2.54</b>	1.87	1.33
(14,10)	<b>2.47</b>	1.89	1.38
(15,10)	<b>2.00</b>	1.48	1.09
(16,10)	<b>1.77</b>	1.30	0.91

Table 10: Decoding throughputs (GB/s): Three parities

$(n, k, w)$	Proposed	Vectorized XOR-based CRS code	STAR code [8]	Quantcast QFS [13]
(8,5,4)	<b>4.28</b>	3.08	3.20	1.77
(8,5,8)	<b>3.04</b>	2.11	3.20	1.77
(9,6,4)	<b>4.13</b>	3.41	3.23	1.74
(9,6,8)	<b>2.94</b>	2.20	3.23	1.74
(10,7,4)	<b>4.55</b>	3.53	3.52	1.77
(10,7,8)	<b>3.28</b>	2.29	3.52	1.77
(11,8,4)	<b>4.70</b>	3.78	3.13	1.68
(11,8,8)	<b>3.08</b>	2.31	3.13	1.68
(13,10,4)	<b>4.86</b>	3.71	3.50	1.71
(13,10,8)	3.21	2.37	<b>3.50</b>	1.71

Table 11: Decoding throughputs (GB/s): Two parities

$(n, k, w)$	Proposed	Vectorized XOR-based CRS code	Vectorized Raid-6	EVEN ODD [2]	RDP [6]
(7,5,4)	5.52	4.85	n/a	6.66	<b>7.28</b>
(7,5,8)	3.87	4.65	2.64	6.66	<b>7.28</b>
(8,6,4)	5.43	5.14	n/a	7.42	<b>8.00</b>
(8,6,8)	5.45	4.86	2.67	7.42	<b>8.00</b>
(9,7,4)	6.03	5.37	n/a	7.65	<b>8.13</b>
(9,7,8)	4.46	5.06	2.73	7.65	<b>8.13</b>
(10,8,4)	5.88	5.73	n/a	7.93	<b>8.44</b>
(10,8,8)	4.89	5.11	2.77	7.93	<b>8.44</b>
(12,10,4)	6.23	5.89	n/a	7.49	<b>9.10</b>
(12,10,8)	4.45	5.52	2.81	7.49	<b>9.10</b>



# Contents

- Motivation
- Background and Review
- Individual Techniques and Combinations
- Bitmatrix Optimization
- Proposed Coding Procedure and Evaluation
- Conclusion



# Conclusion

- Comprehensive study of acceleration techniques
- Combine existing techniques and jointly optimize the bitmatrix
- Proposed approach outperforms most existing approaches in encoding throughput.



# Conclusion : Key Finding

- Vectorization at XOR-level is much better choice than vectorization of finite field operations
  - Higher throughput (~ 2x faster)
  - Easy migration to newer SIMD ISA

SSE: 128 bit

AVX2: 256 bit

AVX-512: 512 bit

`_mm_xor_si128` → `_mm256_xor_si256` → `_mm512_xor_epi32`



# Thank you!

*Questions?*

Tianli Zhou: zhoutl1106@tamu.edu

Chao Tian: chao.tian@tamu.edu

Source code:

<https://bitbucket.org/zhoutl1106/zerasure.git>

<https://github.com/zhoutl1106/zerasure.git>

