



SC21

St. Louis, MO | science & beyond.

Accelerating XOR-Based Erasure Coding using Program Optimization Techniques

Yuya Uezato

DWANGO, Co., Ltd.



“Accelerating Erasure Coding (EC)” means ...

Optimizing matrix multiplication over two finite fields:

$$\text{(for Standard EC)} \quad A \times B \quad \text{over } \mathbb{F}[2^8],$$

$$\text{(for XOR-based EC)} \quad C \times D \quad \text{over } \mathbb{F}[2].$$

- ▶ *Byte* Finite Field $\mathbb{F}[2^8]$ is a field with 256 elements.
- ▶ *Bit* Finite Field $\mathbb{F}[2] = \{0, 1\}$ is a field with the two bits.

“Accelerating Erasure Coding (EC)” means ...

What we need to know about $\mathbb{F}[2^8]$ and $\mathbb{F}[2]$.

Optimizing

$\mathbb{F}[2^8]$ is a field with $2^8 = 256$ elements.

★ 1-byte (8-bits) data can be seen as an element of $\mathbb{F}[2^8]$.

▶ The definition is complex.
(We will see it in the later page).

▶ *Byte*

▶ *Bit*

“Accelerating Erasure Coding (EC)” means ...

What we need to know about $\mathbb{F}[2^8]$ and $\mathbb{F}[2]$.

Optimizing

$\mathbb{F}[2^8]$ is a field with $2^8 = 256$ elements.

★ 1-byte (8-bits) data can be seen as an element of $\mathbb{F}[2^8]$.

▶ The definition is complex.
(We will see it in the later page).

▶ *Byte*

▶ *Bit*

$\mathbb{F}[2] = \{0, 1\}$ is a field of bits.

▶ Its addition is XOR \oplus .

▶ Its multiplication is AND $\&$.

▶ 0 and 1 satisfy the following:

$$x \oplus 0 = 0 \oplus x = x, \quad y \& 1 = 1 \& y = y.$$

“Accelerating Erasure Coding (EC)” means ...

Optimizing matrix multiplication over two finite fields:

$$\text{(for Standard EC)} \quad A \times B \quad \text{over } \mathbb{F}[2^8],$$

$$\text{(for XOR-based EC)} \quad C \times D \quad \text{over } \mathbb{F}[2].$$

- ▶ *Byte* Finite Field $\mathbb{F}[2^8]$ is a field with 256 elements.
- ▶ *Bit* Finite Field $\mathbb{F}[2] = \{0, 1\}$ is a field with the two bits.

A is small. B is large.

“Accelerating Erasure Coding (EC)” means ...

Optimizing matrix multiplication over two finite fields:

$$\text{(for Standard EC)} \quad A \times B \quad \text{over } \mathbb{F}[2^8],$$

$$\text{(for XOR-based EC)} \quad C \times D \quad \text{over } \mathbb{F}[2].$$

- ▶ *Byte* Finite Field $\mathbb{F}[2^8]$ is a field with 256 elements.
- ▶ *Bit* Finite Field $\mathbb{F}[2] = \{0, 1\}$ is a field with the two bits.

A is small. B is large. In this talk, as an example, we consider:

$$\begin{array}{c} 10 \\ \boxed{A} \\ 14 \end{array} \times_{\mathbb{F}[2^8]} \begin{array}{c} 1\text{M}-4\text{M} \\ \boxed{B} \\ 10 \end{array}$$

“Accelerating Erasure Coding (EC)” means ...

Optimizing matrix multiplication over two finite fields:

$$\text{(for Standard EC)} \quad A \times B \quad \text{over } \mathbb{F}[2^8],$$

$$\text{(for XOR-based EC)} \quad C \times D \quad \text{over } \mathbb{F}[2].$$

- ▶ *Byte* Finite Field $\mathbb{F}[2^8]$ is a field with 256 elements.
- ▶ *Bit* Finite Field $\mathbb{F}[2] = \{0, 1\}$ is a field with the two bits.

A is small. B is large. In this talk, as an example, we consider:

$$\begin{array}{c} 10 \\ \boxed{A} \\ 14 \end{array} \times_{\mathbb{F}[2^8]} \begin{array}{c} 1M-4M \\ \boxed{B} \\ 10 \end{array}$$

This setting comes from *erasure coding*.

What are Erasure Coding (EC) and XOR-Based EC

Erasure Coding: Method for Data Redundancy

Example (Building a streaming media server with criteria)

1. We have 14 nodes. Each node has a 20TB disk.
2. We can load data even if nodes ≤ 4 are down.
3. The total capacity of our server = 200TB.
 - ▶ $14 \cdot 20 - 200 = 80\text{TB}$ can be used for data redundancy.

Erasure Coding: Method for Data Redundancy

Example (Building a streaming media server with criteria)

1. We have 14 nodes. Each node has a 20TB disk.
2. We can load data even if nodes ≤ 4 are down.
3. The total capacity of our server = 200TB.
 - ▶ $14 \cdot 20 - 200 = 80\text{TB}$ can be used for data redundancy.

For this criteria, we can employ Reed-Solomon EC $\mathbf{RS}(d = 10, p = 4)$.

- ▶ d : we can assume d -nodes are living. ($d = 14 - p = 10$).
- ▶ p : we can permit nodes $\leq p$ go down. ($p = 4$).

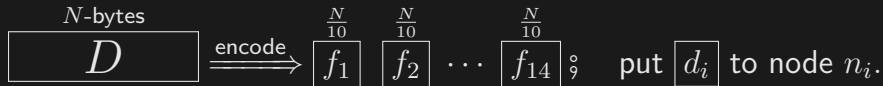
Erasure Coding: Method for Data Redundancy

Example (Building a streaming media server with criteria)

1. We have 14 nodes. Each node has a 20TB disk.
2. We can load data even if nodes ≤ 4 are down.
3. The total capacity of our server = 200TB.
 - ▶ $14 \cdot 20 - 200 = 80\text{TB}$ can be used for data redundancy.

For this criteria, we can employ Reed-Solomon EC $\mathbf{RS}(d = 10, p = 4)$.

- ▶ d : we can assume d -nodes are living. ($d = 14 - p = 10$).
- ▶ p : we can permit nodes $\leq p$ go down. ($p = 4$).



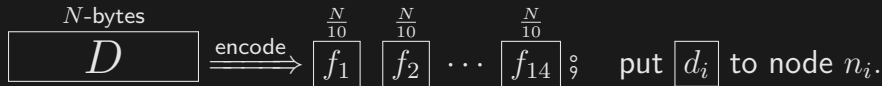
Erasure Coding: Method for Data Redundancy

Example (Building a streaming media server with criteria)

1. We have 14 nodes. Each node has a 20TB disk.
2. We can load data even if nodes ≤ 4 are down.
3. The total capacity of our server = 200TB.
 - ▶ $14 \cdot 20 - 200 = 80\text{TB}$ can be used for data redundancy.

For this criteria, we can employ Reed-Solomon EC $\mathbf{RS}(d = 10, p = 4)$.

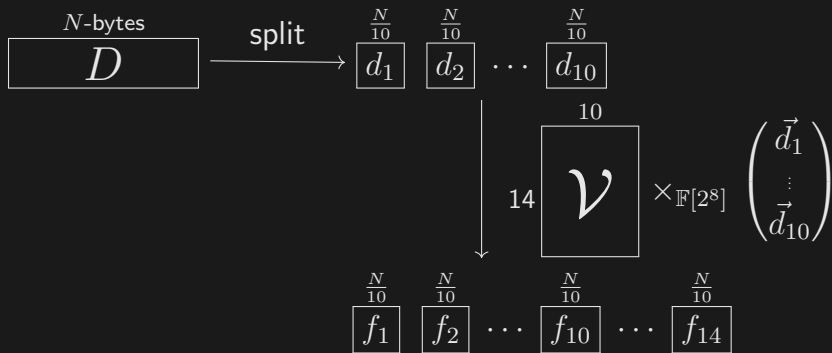
- ▶ d : we can assume d -nodes are living. ($d = 14 - p = 10$).
- ▶ p : we can permit nodes $\leq p$ go down. ($p = 4$).



How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?

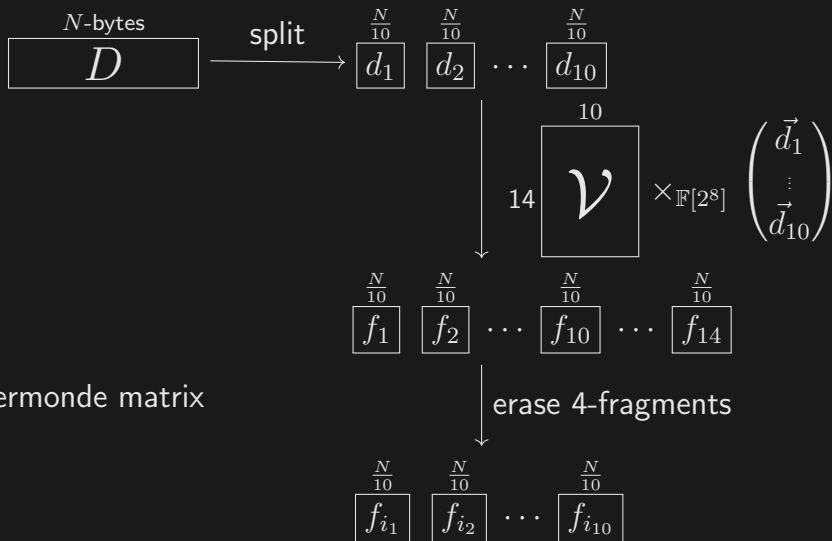


How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?



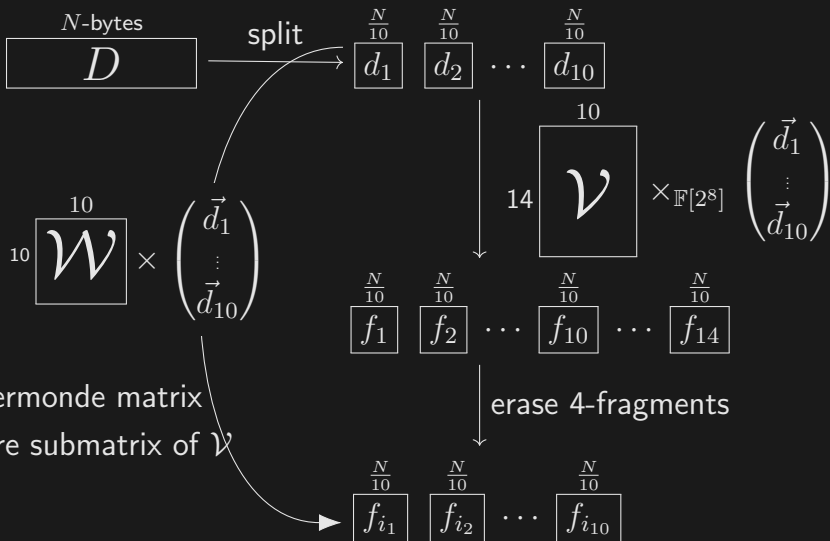
► \mathcal{V} : Vandermonde matrix

How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?



► \mathcal{V} : Vandermonde matrix

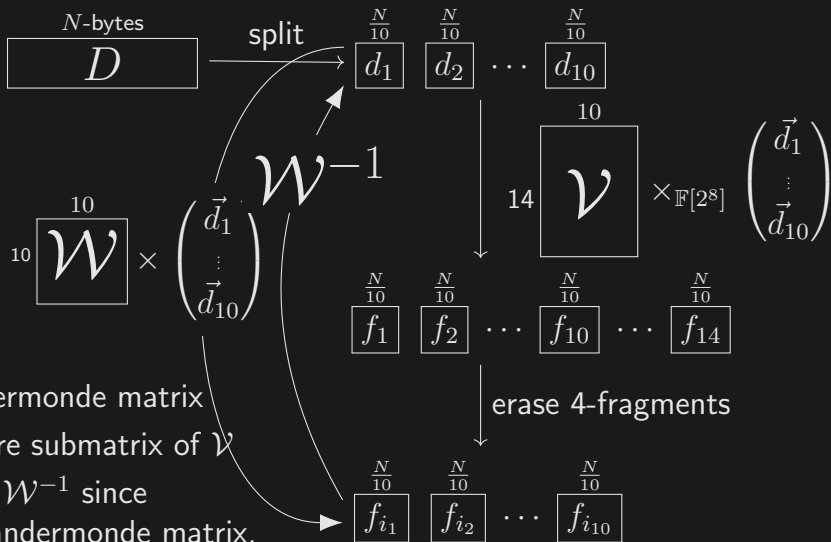
How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?



► \mathcal{V} : Vandermonde matrix

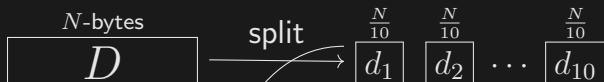
► \mathcal{W} : square submatrix of \mathcal{V}

How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?



- ▶ \mathcal{V} : Vandermonde matrix
- ▶ \mathcal{W} : square submatrix of \mathcal{V}
- ▶ We have \mathcal{W}^{-1} since \mathcal{V} is a Vandermonde matrix.

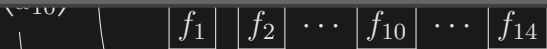
How encoding and decoding are implemented in $\mathbf{RS}(10, 4)$?



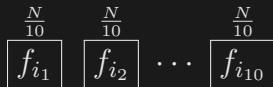
How large is N in a real application?

In my company, D is a short video whose size is 10MB–40MB:

- ▶ The size of 10 secs videos of 1080p & 30fps \sim 12MB.
- ▶ The size of 5 secs videos of 4K & 30fps \sim 35MB.



erase 4-fragments



- ▶ \mathcal{V} : Vandermonde matrix
- ▶ \mathcal{W} : square submatrix of \mathcal{V}
- ▶ We have \mathcal{W}^{-1} since \mathcal{V} is a Vandermonde matrix.

Optimizing $\mathcal{V} \times_{\mathbb{F}[2^8]} D$

Q. What is the heaviest operation on $\mathcal{V} \times_{\mathbb{F}[2^8]} D$?

A. Multiplication of $\mathbb{F}[2^8]$:

Optimizing $\mathcal{V} \times_{\mathbb{F}[2^8]} D$

Q. What is the heaviest operation on $\mathcal{V} \times_{\mathbb{F}[2^8]} D$?

A. Multiplication of $\mathbb{F}[2^8]$:

► Internally, $p \in \mathbb{F}[2^8]$ is a 7-degree polynomial over $\mathbb{F}[2]$:

$$b_7x^7 + b_6x^6 + \cdots + b_1x + b_0 \quad \text{where } b_i \in \mathbb{F}[2].$$

Optimizing $\mathcal{V} \times_{\mathbb{F}[2^8]} D$

Q. What is the heaviest operation on $\mathcal{V} \times_{\mathbb{F}[2^8]} D$?

A. Multiplication of $\mathbb{F}[2^8]$:

- ▶ Internally, $p \in \mathbb{F}[2^8]$ is a 7-degree polynomial over $\mathbb{F}[2]$:

$$b_7x^7 + b_6x^6 + \cdots + b_1x + b_0 \quad \text{where } b_i \in \mathbb{F}[2].$$

- ▶ $p_1 + p_2$ of $\mathbb{F}[2^8]$ is the polynomial addition.
Easy because just componentwise XOR:

$$(b_7 \oplus b'_7)x^7 + (b_6 \oplus b'_6)x^6 + \cdots + (b_0 \oplus b'_0).$$

Optimizing $\mathcal{V} \times_{\mathbb{F}[2^8]} D$

Q. What is the heaviest operation on $\mathcal{V} \times_{\mathbb{F}[2^8]} D$?

A. Multiplication of $\mathbb{F}[2^8]$:

- ▶ Internally, $p \in \mathbb{F}[2^8]$ is a 7-degree polynomial over $\mathbb{F}[2]$:

$$b_7x^7 + b_6x^6 + \dots + b_1x + b_0 \quad \text{where } b_i \in \mathbb{F}[2].$$

- ▶ $p_1 + p_2$ of $\mathbb{F}[2^8]$ is the polynomial addition.
Easy because just componentwise XOR:

$$(b_7 \oplus b'_7)x^7 + (b_6 \oplus b'_6)x^6 + \dots + (b_0 \oplus b'_0).$$

- ▶ On the other hand, $p_1 \cdot p_2$ of $\mathbb{F}[2^8]$ is CPU-heavy and slow:
 1. We do the 7-degree polynomial multiplication $p_1 \times p_2$.
 2. We take the modulo by a special polynomial $(p_1 \times p_2) \bmod p$.

Optimizing $\mathcal{V} \times_{\mathbb{F}[2^8]} D$

Q. What is the heaviest operation on $\mathcal{V} \times_{\mathbb{F}[2^8]} D$?

A. Multiplication of $\mathbb{F}[2^8]$:

- ▶ Internally, $p \in \mathbb{F}[2^8]$ is a 7-degree polynomial over $\mathbb{F}[2]$:

$$b_7x^7 + b_6x^6 + \dots + b_1x + b_0 \quad \text{where } b_i \in \mathbb{F}[2].$$

- ▶ $p_1 + p_2$ of $\mathbb{F}[2^8]$ is the polynomial addition.

Easy because just componentwise XOR:

$$(b_7 \oplus b'_7)x^7 + (b_6 \oplus b'_6)x^6 + \dots + (b_0 \oplus b'_0).$$

- ▶ On the other hand, $p_1 \cdot p_2$ of $\mathbb{F}[2^8]$ is CPU-heavy and slow:

1. We do the 7-degree polynomial multiplication $p_1 \times p_2$.
2. We take the modulo by a special polynomial $(p_1 \times p_2) \bmod p$.

XOR-based EC is one way to vanish \cdot of $\mathbb{F}[2^8]$.

XOR-based EC: From $\mathbb{F}[2^8]$ to BitMatrix ($\mathbb{F}[2]$ -Matrix)

- ▶ 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \mathbf{8} \overset{1}{\square} \mathbb{F}[2]$.

XOR-based EC: From $\mathbb{F}[2^8]$ to BitMatrix ($\mathbb{F}[2]$ -Matrix)

- ▶ 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \begin{matrix} 1 \\ \boxed{\mathbb{F}[2]} \\ 8 \end{matrix}$.
- ▶ There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \begin{matrix} 8 \\ \boxed{\mathbb{F}[2]} \\ 8 \end{matrix}$ i.e.,

XOR-based EC: From $\mathbb{F}[2^8]$ to BitMatrix ($\mathbb{F}[2]$ -Matrix)

- ▶ 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \mathfrak{s} \begin{matrix} 1 \\ \boxed{\mathbb{F}[2]} \end{matrix}$.
- ▶ There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \mathfrak{s} \begin{matrix} 8 \\ \boxed{\mathbb{F}[2]} \end{matrix}$ i.e.,

$$\forall x, y \in \mathbb{F}[2^8]. \begin{cases} x + y = \mathcal{B}^{-1}(\mathcal{B}(x) + \mathcal{B}(y)), \\ x \cdot y = \mathcal{B}^{-1}(\mathcal{B}(x) \times \mathcal{B}(y)) \end{cases}$$

XOR-based EC: From $\mathbb{F}[2^8]$ to BitMatrix ($\mathbb{F}[2]$ -Matrix)

- ▶ 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \mathfrak{s} \begin{matrix} 1 \\ \boxed{\mathbb{F}[2]} \end{matrix}$.
- ▶ There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \mathfrak{s} \begin{matrix} 8 \\ \boxed{\mathbb{F}[2]} \end{matrix}$ i.e.,

$$\forall x, y \in \mathbb{F}[2^8]. \begin{cases} x + y = \mathcal{B}^{-1}(\mathcal{B}(x) + \mathcal{B}(y)), \\ x \cdot y = \mathcal{B}^{-1}(\mathcal{B}(x) \times \mathcal{B}(y)) \end{cases}$$

Prop: Emulate $\mathcal{W}^{-1} \times (\mathcal{W} \times D) = D$ in the $\mathbb{F}[2]$ world

$$\mathcal{B}(\mathcal{W}^{-1}) \begin{matrix} \mathbb{F}[2] \\ \times \end{matrix} (\mathcal{B}(\mathcal{W}) \begin{matrix} \mathbb{F}[2] \\ \times \end{matrix} \tilde{D}) = \mathcal{B}(\mathcal{W}^{-1} \times \mathcal{W}) \times \tilde{D} = \tilde{D}.$$

XOR-based EC: From $\mathbb{F}[2^8]$ to BitMatrix ($\mathbb{F}[2]$ -Matrix)

► 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \mathfrak{s} \boxed{\mathbb{F}[2]}^1$.

► There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \mathfrak{s} \boxed{\mathbb{F}[2]}^8$

$$\begin{pmatrix} x_1 & x_2 \\ x_3 & x_4 \end{pmatrix} \times_{\mathbb{F}[2^8]} \begin{pmatrix} d_1 & \cdots \\ d_2 & \cdots \end{pmatrix} = \begin{pmatrix} x_1 \cdot d_1 + x_2 \cdot d_2 & \cdots \\ x_3 \cdot d_1 + x_4 \cdot d_4 & \cdots \end{pmatrix}$$

⇓

$$\begin{pmatrix} 1 & 1 & 0 & 1 & \cdots \\ 0 & 0 & 1 & 1 & \cdots \\ 0 & 1 & 1 & 0 & \cdots \\ 1 & 0 & 0 & 1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \\ \vec{x}_4 \\ \vdots \end{pmatrix} = \begin{pmatrix} \vec{x}_1 \oplus \vec{x}_2 \oplus \vec{x}_4 \oplus \cdots \\ \vec{x}_3 \oplus \vec{x}_4 \oplus \cdots \\ \vec{x}_2 \oplus \vec{x}_3 \oplus \cdots \\ \vec{x}_1 \oplus \vec{x}_4 \oplus \cdots \\ \vdots \end{pmatrix}$$

\oplus is byte-array XOR.



Comparing MM over $\mathbb{F}[2^8]$ and MM over $\mathbb{F}[2]$ for Encoding


Trade-off in Matrix Multiplication	RS(10, 4) by $\mathbb{F}[2^8]$	RS(10, 4) by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \overset{10}{\boxed{\mathbb{F}[2^8]}}$	$\mathcal{B}(\mathcal{V}) : 112 \overset{80}{\boxed{\mathbb{F}[2]}}$
Speed of Core Operation	+ of $\mathbb{F}[2^8]$ is fast · of $\mathbb{F}[2^8]$ is slow	bytevec-XOR \oplus is fast (SIMDable)


Comparing MM over $\mathbb{F}[2^8]$ and MM over $\mathbb{F}[2]$ for Encoding

Trade-off in Matrix Multiplication	$\mathbf{RS}(10, 4)$ by $\mathbb{F}[2^8]$	$\mathbf{RS}(10, 4)$ by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \overset{10}{\boxed{\mathbb{F}[2^8]}}$	$\mathcal{B}(\mathcal{V}) : 112 \overset{80}{\boxed{\mathbb{F}[2]}}$
Speed of Core Operation	+ of $\mathbb{F}[2^8]$ is fast · of $\mathbb{F}[2^8]$ is slow	bytevec-XOR \oplus is fast (SIMDable)

Encoding Throughput Comparison (on Intel CPU):

GB/s	ISA-L  $\mathbb{F}[2^8]$	State-of-the-art  $\mathbb{F}[2]$
RS(10, 4)	6.79	4.94
RS(10, 3)	6.78	6.15
RS(9, 3)	7.31	6.17



 ISA-L: Intel's EC library <https://github.com/intel/isa-1>


 T. Zhou & C. Tian. 2020. *Fast Erasure Coding for Data Storage: A Comprehensive Study of the Acceleration Techniques.*


Comparing MM over $\mathbb{F}[2^8]$ and MM over $\mathbb{F}[2]$ for Encoding

Trade-off in Matrix Multiplication	$\mathbf{RS}(10, 4)$ by $\mathbb{F}[2^8]$	$\mathbf{RS}(10, 4)$ by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \overset{10}{\boxed{\mathbb{F}[2^8]}}$	$\mathcal{B}(\mathcal{V}) : 112 \overset{80}{\boxed{\mathbb{F}[2]}}$
Speed of Core Operation	+ of $\mathbb{F}[2^8]$ is fast · of $\mathbb{F}[2^8]$ is slow	bytevec-XOR \oplus is fast (SIMDable)

Encoding Throughput Comparison (on Intel CPU):

GB/s	ISA-L  $\mathbb{F}[2^8]$	State-of-the-art  $\mathbb{F}[2]$	Ours(New!) $\mathbb{F}[2]$
RS(10, 4)	6.79	4.94	8.92
RS(10, 3)	6.78	6.15	11.78
RS(9, 3)	7.31	6.17	11.97

 ISA-L: Intel's EC library <https://github.com/intel/isa-1>

 T. Zhou & C. Tian. 2020. *Fast Erasure Coding for Data Storage: A Comprehensive Study of the Acceleration Techniques.*

Our Contribution:

Optimizing Bitmatrix Multiplication

as

Program Optimization Problem

MM over $\mathbb{F}[2]$ = Running Straight Line Program

We identify bitmatrix multiplication as *straight line program* (SLP):

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix} \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{a} \\ \vec{b} \\ \vec{c} \\ \vec{d} \end{pmatrix} \Downarrow$$

$P(a, b, c, d)$

$v_1 \leftarrow a \oplus b;$

$v_2 \leftarrow a \oplus b \oplus c;$

$v_3 \leftarrow b \oplus c \oplus d;$

return(v_1, v_2, v_3)

MM over $\mathbb{F}[2]$ = Running Straight Line Program

We identify bitmatrix multiplication as *straight line program* (SLP):

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix} \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{a} \\ \vec{b} \\ \vec{c} \\ \vec{d} \end{pmatrix} = \begin{pmatrix} \vec{a} \oplus \vec{b} \\ \vec{a} \oplus \vec{b} \oplus \vec{c} \\ \vec{b} \oplus \vec{c} \oplus \vec{d} \end{pmatrix}$$

\Downarrow

$$\frac{P(a, b, c, d)}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow b \oplus c \oplus d; \\ \text{return}(v_1, v_2, v_3)}{[[P]] = \text{return}(v_1, v_2, v_3) \\ = \langle a \oplus b, \\ a \oplus b \oplus c, \\ b \oplus c \oplus d \rangle}$$

MM over $\mathbb{F}[2]$ = Running Straight Line Program

We identify bitmatrix multiplication as *straight line program* (SLP):

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix} \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{a} \\ \vec{b} \\ \vec{c} \\ \vec{d} \end{pmatrix} = \begin{pmatrix} \vec{a} \oplus \vec{b} \\ \vec{a} \oplus \vec{b} \oplus \vec{c} \\ \vec{b} \oplus \vec{c} \oplus \vec{d} \end{pmatrix}$$

\Downarrow

$$\frac{P(a, b, c, d)}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow b \oplus c \oplus d; \\ \text{return}(v_1, v_2, v_3)}{\llbracket P \rrbracket = \text{return}(v_1, v_2, v_3) \\ = \langle a \oplus b, \\ a \oplus b \oplus c, \\ b \oplus c \oplus d \rangle}$$

★ "Bitmatrix as SLP" is not a new idea (See. Boyar+ 2008)

MM over $\mathbb{F}[2]$ = Running Straight Line Program

We identify bitmatrix multiplication as *straight line program* (SLP):

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix} \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{a} \\ \vec{b} \\ \vec{c} \\ \vec{d} \end{pmatrix} = \begin{pmatrix} \vec{a} \oplus \vec{b} \\ \vec{a} \oplus \vec{b} \oplus \vec{c} \\ \vec{b} \oplus \vec{c} \oplus \vec{d} \end{pmatrix}$$

$$\frac{P(a, b, c, d)}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow b \oplus c \oplus d; \\ \text{return}(v_1, v_2, v_3)}{[[P]] = \text{return}(v_1, v_2, v_3) \\ = \langle a \oplus b, \\ a \oplus b \oplus c, \\ b \oplus c \oplus d \rangle}$$

- ★ "Bitmatrix as SLP" is not a new idea (See. Boyar+ 2008)
- ▶ SLP only allow assignments with one kind *binary* operator \oplus .
- ▶ SLP do not have functions, if-branchings, and while-loop, etc.

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \underline{P \quad \#_{\oplus} = 8} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array} \quad \Rightarrow \quad \underline{Q}$$

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \frac{P \quad \#_{\oplus} = 8}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d;} \quad \Rightarrow \quad \frac{Q}{v_1 \leftarrow a \oplus b;} \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array}$$

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \frac{P \quad \#_{\oplus} = 8}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d;} \quad \Rightarrow \quad \frac{Q}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow v_1 \oplus c;} \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array}$$

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \overline{P \quad \#_{\oplus} = 8} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \quad \Rightarrow \quad \begin{array}{l} \overline{Q} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow v_1 \oplus c; \\ v_3 \leftarrow v_2 \oplus d; \end{array}$$

$\text{return}(v_1, v_2, v_3, v_4)$

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \text{P} \quad \#_{\oplus} = 8 \\ \hline v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array} \quad \Rightarrow \quad \begin{array}{l} \text{Q} \quad \#_{\oplus} = 4 \\ \hline v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow v_1 \oplus c; \\ v_3 \leftarrow v_2 \oplus d; \\ v_4 \leftarrow v_3 \oplus a; \\ \\ \because (a \oplus b \oplus c \oplus d) \oplus a = b \oplus c \oplus d. \\ \text{return}(v_1, v_2, v_3, v_4) \end{array}$$

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \overline{P \quad \#_{\oplus} = 8} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array} \quad \Rightarrow \quad \begin{array}{l} \overline{Q \quad \#_{\oplus} = 4} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow v_1 \oplus c; \\ v_3 \leftarrow v_2 \oplus d; \\ v_4 \leftarrow v_3 \oplus a; \\ \\ \because (a \oplus b \oplus c \oplus d) \oplus a = b \oplus c \oplus d. \\ \text{return}(v_1, v_2, v_3, v_4) \end{array}$$

- ▶ P and Q are equivalent: $\llbracket P \rrbracket = \llbracket Q \rrbracket$.
- ▶ Intuitively, Q ($\#_{\oplus}(Q) = 4$) runs faster than P ($\#_{\oplus}(P) = 8$).

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$\begin{array}{l} \overline{P \quad \#_{\oplus} = 8} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \\ \text{return}(v_1, v_2, v_3, v_4) \end{array} \quad \Rightarrow \quad \begin{array}{l} \overline{Q \quad \#_{\oplus} = 4} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow v_1 \oplus c; \\ v_3 \leftarrow v_2 \oplus d; \\ v_4 \leftarrow v_3 \oplus a; \\ \\ \because (a \oplus b \oplus c \oplus d) \oplus a = b \oplus c \oplus d. \\ \text{return}(v_1, v_2, v_3, v_4) \end{array}$$

- ▶ P and Q are equivalent: $\llbracket P \rrbracket = \llbracket Q \rrbracket$.
- ▶ Intuitively, Q ($\#_{\oplus}(Q) = 4$) runs faster than P ($\#_{\oplus}(P) = 8$).

*Question. For a given SLP P ,
can we quickly find the most efficient equivalent SLP Q ?*

XOR Optimization: Reducing XORs

Optimization Metric $\#_{\oplus}(-)$: the number of XORs.

$$P \quad \#_{\oplus} = 8$$

$$v_1 \leftarrow a \oplus b;$$

$$v_2 \leftarrow a \oplus b \oplus c;$$

$$v_3 \leftarrow a \oplus b \oplus c \oplus d;$$

$$v_4 \leftarrow b \oplus c \oplus d;$$

return(v_1, v_2, v_3, v_4)



$$Q \quad \#_{\oplus} = 4$$

$$v_1 \leftarrow a \oplus b;$$

$$v_2 \leftarrow v_1 \oplus c;$$

$$v_3 \leftarrow v_2 \oplus d;$$

$$v_4 \leftarrow v_3 \oplus a;$$

$$\because (a \oplus b \oplus c \oplus d) \oplus a = b \oplus c \oplus d.$$

return(v_1, v_2, v_3, v_4)

- ▶ P and Q are equivalent: $\llbracket P \rrbracket = \llbracket Q \rrbracket$.
- ▶ Intuitively, Q ($\#_{\oplus}(Q) = 4$) runs faster than P ($\#_{\oplus}(P) = 8$).

Theorem (Boyar+ 2013)

Unless $\mathbf{P} = \mathbf{NP}$, for a given SLP P , in polynomial time, we cannot find Q such that $\llbracket P \rrbracket = \llbracket Q \rrbracket$ and minimizes $\#_{\oplus}(Q)$.

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

REPAIR = **Repeat** PAIR. The key operation is PAIR:

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 8 \end{array} \xrightarrow{\text{PAIR}(a, c)} \begin{array}{l} t_1 \leftarrow a \oplus c; \\ \hline v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus b; \\ v_3 \leftarrow t_1 \oplus b \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 7 \end{array}$$

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

REPAIR = **R**epat PAIR. The key operation is PAIR:

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 8 \end{array} \xrightarrow{\text{PAIR}(a, c)} \begin{array}{l} \overline{t_1 \leftarrow a \oplus c;} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus b; \\ v_3 \leftarrow t_1 \oplus b \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 7 \end{array}$$

How do we choose a pair of terms to do pairing?

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

REPAIR = **R**epeat PAIR. The key operation is PAIR:

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 8 \end{array} \xrightarrow{\text{PAIR}(a, c)} \begin{array}{l} t_1 \leftarrow a \oplus c; \\ \hline v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus b; \\ v_3 \leftarrow t_1 \oplus b \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 7 \end{array}$$

How do we choose a pair of terms to do pairing? *Greedy*.

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \xrightarrow{\text{PAIR}(a, b)} \begin{array}{l} t_1 \leftarrow a \oplus b; \\ \hline v_2 \leftarrow t_1 \oplus c; \\ v_3 \leftarrow t_1 \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 6 \end{array}$$

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

REPAIR = **R**epeat PAIR. The key operation is PAIR:

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 8 \end{array} \xrightarrow{\text{PAIR}(a, c)} \begin{array}{l} t_1 \leftarrow a \oplus c; \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus b; \\ v_3 \leftarrow t_1 \oplus b \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \quad \#_{\oplus} = 7$$

$$\begin{array}{l} t_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus c; \\ v_3 \leftarrow t_1 \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \xrightarrow{\text{PAIR}(t_1, c)} \begin{array}{l} t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ v_3 \leftarrow t_2 \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \xrightarrow{\text{PAIR}(b, c)} \begin{array}{l} t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ t_3 \leftarrow b \oplus c; \\ v_3 \leftarrow t_2 \oplus d; \\ v_4 \leftarrow t_3 \oplus d; \end{array}$$

Our Heuristic: Grammar Compression Algorithm REPAIR

Originally, REPAIR is an algorithm to compress context-free grammars. We use it identifying SLPs as commutative CFGs.

- ▶ Larsson & Moffat. 1999. *Offline dictionary-based compression*
- ▶ Paar. 1997. *Optimized arithmetic for Reed-Solomon encoders*

REPAIR = Repeat PAIR. The key operation is PAIR:

$$\begin{array}{l} v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow a \oplus b \oplus c \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \\ \#_{\oplus} = 8 \end{array} \xrightarrow{\text{PAIR}(a, c)} \begin{array}{l} \overline{t_1 \leftarrow a \oplus c;} \\ v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow t_1 \oplus b; \\ v_3 \leftarrow t_1 \oplus b \oplus d; \\ v_4 \leftarrow b \oplus c \oplus d; \end{array} \quad \#_{\oplus} = 7$$

The commutative version of REPAIR accommodates

Commutativity : $x \oplus y = y \oplus x$, Associativity : $(x \oplus y) \oplus z = x \oplus (y \oplus z)$.

In the paper, we extend it to XORREPAIR by accommodating

Cancellativity: $x \oplus x \oplus y = y$.

Memory Access Optimization: MultiSLP

Optimization Metric: $\#_{\text{mem}}(-)$ = the number of memory access.

Quiz: How many times will this program access memory?

$$\#_{\text{mem}} \left(v \leftarrow A \oplus B \oplus C \oplus D \right) = ?$$

Memory Access Optimization: MultiSLP

Optimization Metric: $\#_{\text{mem}}(-)$ = the number of memory access.

Quiz: How many times will this program access memory?

$$\#_{\text{mem}} \left(v \leftarrow A \oplus B \oplus C \oplus D \right) = 9$$

because each \oplus issues two read and one write:

$$t_1 \leftarrow A \oplus B; \quad t_2 \leftarrow t_1 \oplus C; \quad v \leftarrow t_2 \oplus D;$$

Memory Access Optimization: **MultiSLP**

Optimization Metric: $\#_{\text{mem}}(-)$ = the number of memory access.

Quiz: How many times will this program access memory?

$$\#_{\text{mem}} \left(v \leftarrow A \oplus B \oplus C \oplus D \right) = 9$$

because each \oplus issues two read and one write:

$$t_1 \leftarrow A \oplus B; \quad t_2 \leftarrow t_1 \oplus C; \quad v \leftarrow t_2 \oplus D;$$

t_1 and t_2 are wasteful: they are released immediately after allocated.

To reduce such wastefulness,
we extend SLP to *MultiSLP*, which allows n -arity XORs.

Memory Access Optimization: MultiSLP

Optimization Metric: $\#_{\text{mem}}(-)$ = the number of memory access.

Quiz: How many times will this program access memory?

$$\#_{\text{mem}} \left(v \leftarrow A \oplus B \oplus C \oplus D \right) = 9$$

because each \oplus issues two read and one write:

$$t_1 \leftarrow A \oplus B; \quad t_2 \leftarrow t_1 \oplus C; \quad v \leftarrow t_2 \oplus D;$$

On MultiSLP, we can

$$v \leftarrow \oplus_4(A, B, C, D);$$

Thus, we have $\#_{\text{mem}} = 5$.

```
 $\oplus_4(A, B, C, D: [\text{byte}]) \{$   
  var  $v = \text{Array}::\text{new}(A.\text{len})$ ;  
  for  $i$  in  $[0..A.\text{len})$ :  
    byte  $r = A[i] \wedge B[i]$   
     $r = r \wedge C[i]$ ;  
     $v[i] = r \wedge D[i]$ ;  
  return  $v$ ;  
}
```

New Metric and Memory Optimization Problem

From a given P , can we quickly (= in polynomial time)
find an equivalent and most memory efficient Q w.r.t. $\#_{\text{mem}}$?

$$P : \begin{aligned} v_1 &\leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 &\leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \#_{\text{mem}}(P) &= 24 \end{aligned}$$

New Metric and Memory Optimization Problem

From a given P , can we quickly (= in polynomial time)
find an equivalent and most memory efficient Q w.r.t. $\#_{\text{mem}}$?

$$\begin{array}{l} P : \\ v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \#_{\text{mem}}(P) = 24 \end{array} \quad \Longrightarrow \quad \begin{array}{l} Q : \\ t \leftarrow \oplus_4(a, b, c, d); \\ v_1 \leftarrow t \oplus e; \\ v_2 \leftarrow t \oplus f; \\ \#_{\text{mem}}(Q) = 11 \end{array}$$

New Metric and Memory Optimization Problem

From a given P , can we quickly (= in polynomial time) find an equivalent and most memory efficient Q w.r.t. $\#_{\text{mem}}$?

$$\begin{array}{l} P : \quad v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ \quad \quad v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \quad \quad \#_{\text{mem}}(P) = 24 \end{array} \quad \Longrightarrow \quad \begin{array}{l} Q : \quad t \leftarrow \oplus_4(a, b, c, d); \\ \quad \quad v_1 \leftarrow t \oplus e; \\ \quad \quad v_2 \leftarrow t \oplus f; \\ \quad \quad \#_{\text{mem}}(Q) = 11 \end{array}$$

Unfortunately, we showed the following intractability result:

Theorem (Our NEW theoretical result)

Unless $\mathbf{P} = \mathbf{NP}$, for a given SLP P , in polynomial time, we cannot find Q that $\llbracket P \rrbracket = \llbracket Q \rrbracket$ and minimizes $\#_{\text{mem}}(Q)$.

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

$$\alpha \leftarrow \oplus(x_1, \dots, x_n);$$

\vdots

$$\beta \leftarrow \oplus(y_1, \dots, \alpha, \dots, y_m) \xrightarrow{\text{fuse}} \beta \leftarrow \oplus(y_1, \dots, x_1, \dots, x_n, \dots, y_m)$$

☆ α appears once in the program

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

$$\alpha \leftarrow \oplus(x_1, \dots, x_n);$$

\vdots

$$\beta \leftarrow \oplus(y_1, \dots, \alpha, \dots, y_m) \xrightarrow{\text{fuse}} \beta \leftarrow \oplus(y_1, \dots, x_1, \dots, x_n, \dots, y_m)$$

☆ α appears once in the program

Example.

$$v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e;$$

$$v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f;$$

$\#_{\text{mem}}(24)$

$\xrightarrow{\text{REPAIR}}$

$$t_1 \leftarrow a \oplus b;$$

$$t_2 \leftarrow t_1 \oplus c;$$

$$t_3 \leftarrow t_2 \oplus d;$$

$$v_1 \leftarrow t_3 \oplus e;$$

$$v_2 \leftarrow t_3 \oplus f;$$

$\#_{\text{mem}}(15)$

$\xrightarrow{\text{fuse}(t_1)}$

$$t_2 \leftarrow \oplus_3(a, b, c);$$

$$t_3 \leftarrow t_2 \oplus d;$$

$$v_1 \leftarrow t_3 \oplus e;$$

$$v_2 \leftarrow t_3 \oplus f;$$

$\#_{\text{mem}}(13)$

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

$$\begin{array}{l} \alpha \leftarrow \oplus(x_1, \dots, x_n); \\ \vdots \\ \beta \leftarrow \oplus(y_1, \dots, \alpha, \dots, y_m) \end{array} \xrightarrow{\text{fuse}} \beta \leftarrow \oplus(y_1, \dots, x_1, \dots, x_n, \dots, y_m)$$

☆ α appears once in the program

Example.

$$\begin{array}{l} v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \end{array} \xrightarrow{\text{REPAIR}} \begin{array}{l} t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array} \xrightarrow{\text{fuse}(t_1)} \begin{array}{l} t_2 \leftarrow \oplus_3(a, b, c); \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array}$$

#mem(24)

$$\begin{array}{l} t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array} \xrightarrow{\text{fuse}(t_1)} \begin{array}{l} t_2 \leftarrow \oplus_3(a, b, c); \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array}$$

$$\begin{array}{l} t_2 \leftarrow \oplus_3(a, b, c); \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array} \quad \#_{\text{mem}}(13)$$

$$\xrightarrow{\text{fuse}(t_2)} \begin{array}{l} t_3 \leftarrow \oplus_4(a, b, c, d); \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \end{array} \quad \#_{\text{mem}}(11)$$

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

$$\begin{array}{l} \alpha \leftarrow \oplus(x_1, \dots, x_n); \\ \vdots \\ \beta \leftarrow \oplus(y_1, \dots, \alpha, \dots, y_m) \end{array} \xrightarrow{\text{fuse}} \beta \leftarrow \oplus(y_1, \dots, x_1, \dots, x_n, \dots, y_m)$$

☆ α appears once in the program

Example.

$$\begin{array}{l} v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \#_{\text{mem}}(24) \end{array} \xrightarrow{\text{REPAIR}} \begin{array}{l} t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \\ \#_{\text{mem}}(15) \end{array} \xrightarrow{\text{fuse}(t_1)} \begin{array}{l} t_2 \leftarrow \oplus_3(a, b, c); \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \\ \#_{\text{mem}}(13) \end{array}$$
$$\begin{array}{l} t_3 \leftarrow \oplus_4(a, b, c, d); \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \\ \#_{\text{mem}}(11) \end{array} \xrightarrow{\text{fuse}(t_2)} \begin{array}{l} \text{NOT fuse}(t_3) \text{ by } \star \\ v_1 \leftarrow \oplus_5(a, b, c, d, e); \\ v_2 \leftarrow \oplus_5(a, b, c, d, f); \\ \#_{\text{mem}}(12) \end{array}$$

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Example: Calculate $\#_{I/O}(4, P)$ for the following example SLP P :

$v_1 \leftarrow A \oplus B;$ $*_1 *_2 *_3 *_4$

$v_2 \leftarrow \oplus(E, D, A);$

$v_3 \leftarrow v_1 \oplus E;$

$v_4 \leftarrow v_1 \oplus C;$

$\text{return}(v_2, v_3, v_4);$

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Example: Calculate $\#_{I/O}(4, P)$ for the following example SLP P :

$$v_1 \leftarrow A \oplus B; \quad *_{1} *_{2} *_{3} *_{4} \xrightarrow[2]{\mathcal{R}(A)} *_{2} *_{3} *_{4} A$$

$$v_2 \leftarrow \oplus(E, D, A);$$

$$v_3 \leftarrow v_1 \oplus E;$$

$$v_4 \leftarrow v_1 \oplus C;$$

return(v_2, v_3, v_4);

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Example: Calculate $\#_{I/O}(4, P)$ for the following example SLP P :

$v_1 \leftarrow A \oplus B;$ $*_1 * _2 * _3 * _4 \xrightarrow[2]{\mathcal{R}(A)} * _2 * _3 * _4 A \xrightarrow[2]{\mathcal{R}(B)} * _3 * _4 AB$

$v_2 \leftarrow \oplus(E, D, A);$

$v_3 \leftarrow v_1 \oplus E;$

$v_4 \leftarrow v_1 \oplus C;$

return(v_2, v_3, v_4);

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Example: Calculate $\#_{I/O}(4, P)$ for the following example SLP P :

$$v_1 \leftarrow A \oplus B; \quad *_{1} *_{2} *_{3} *_{4} \xrightarrow[2]{\mathcal{R}(A)} *_{2} *_{3} *_{4} A \xrightarrow[2]{\mathcal{R}(B)} *_{3} *_{4} AB \xrightarrow[1]{\mathcal{W}(v_1)}$$

$$v_2 \leftarrow \oplus(E, D, A); \quad *_{4} ABv_1$$

$$v_3 \leftarrow v_1 \oplus E;$$

$$v_4 \leftarrow v_1 \oplus C;$$

return(v_2, v_3, v_4);

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, -)$: the total number of I/O transfers between memory and cache of K -capacity.

We have three kinds of operations for cache:

- ▶ $\mathcal{H}(x)$: Cache Hit for an element x . $\#_{I/O} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{I/O} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{I/O} = 1$.

Example: Calculate $\#_{I/O}(4, P)$ for the following example SLP P :

$$\begin{array}{l}
 v_1 \leftarrow A \oplus B; \quad *_{1} *_{2} *_{3} *_{4} \xrightarrow[2]{\mathcal{R}(A)} *_{2} *_{3} *_{4} A \xrightarrow[2]{\mathcal{R}(B)} *_{3} *_{4} AB \xrightarrow[1]{\mathcal{W}(v_1)} \\
 v_2 \leftarrow \oplus(E, D, A); \quad *_{4} ABv_1 \xrightarrow[2]{\mathcal{R}(E)} ABv_1E \xrightarrow[2]{\mathcal{R}(D)} Bv_1ED \xrightarrow[2]{\mathcal{R}(A)} v_1EDA \xrightarrow[1]{\mathcal{W}(v_2)} \\
 v_3 \leftarrow v_1 \oplus E; \quad EDAv_2 \xrightarrow[2]{\mathcal{R}(v_1)} DAv_2v_1 \xrightarrow[2]{\mathcal{R}(E)} Av_2v_1E \xrightarrow[1]{\mathcal{W}(v_3)} \\
 v_4 \leftarrow v_1 \oplus C; \quad v_2v_1Ev_3 \xrightarrow[0]{\mathcal{H}(v_1)} v_2Ev_3v_1 \xrightarrow[2]{\mathcal{R}(C)} Ev_3v_1C \xrightarrow[1]{\mathcal{W}(v_4)} \\
 \text{return}(v_2, v_3, v_4); \quad v_3v_1Cv_4 \implies \#_{I/O}(4, P) = 20.
 \end{array}$$

First approach: Register Assignment

Idea: Reducing the number of variables can relax the pressure of cache, and thus may reduce $\#_{I/O}$.

We do Recycling variables by *Register assignment*.

First approach: Register Assignment

Idea: Reducing the number of variables can relax the pressure of cache, and thus may reduce $\#_{I/O}$.

We do Recycling variables by *Register assignment*.

$v_1 \leftarrow A \oplus B;$	$\#_{I/O}$ [5]	Register assignment \Longrightarrow	$v_1 \leftarrow A \oplus B;$	$\#_{I/O}$ [5]
$v_2 \leftarrow \oplus(E, D, A);$	[7]		$v_2 \leftarrow \oplus(E, D, A);$	[7]
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$	[5]
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_1 \leftarrow v_1 \oplus C;$	[2]
$\text{return}(v_2, v_3, v_4);$			$\text{return}(v_2, v_3, v_1);$	

First approach: Register Assignment

Idea: Reducing the number of variables can relax the pressure of cache, and thus may reduce $\#_{I/O}$.

We do Recycling variables by *Register assignment*.

	$\#_{I/O}$		$\#_{I/O}$	
$v_1 \leftarrow A \oplus B;$	[5]	Register assignment \Longrightarrow	$v_1 \leftarrow A \oplus B;$	[5]
$v_2 \leftarrow \oplus(E, D, A);$	[7]		$v_2 \leftarrow \oplus(E, D, A);$	[7]
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$	[5]
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_1 \leftarrow v_1 \oplus C;$	[2]
$\text{return}(v_2, v_3, v_4);$			$\text{return}(v_2, v_3, v_1);$	

$$\begin{array}{c}
 \xrightarrow[0]{\mathcal{R}(v_1)} v_2 E v_3 v_1 \xrightarrow[2]{\mathcal{R}(C)} E v_3 v_1 C \xrightarrow[1]{\mathcal{W}(v_4)} v_3 v_1 C v_4 \\
 \Downarrow \\
 \xrightarrow[0]{\mathcal{R}(v_1)} v_2 E v_3 v_1 \xrightarrow[2]{\mathcal{R}(C)} E v_3 v_1 C \xrightarrow[0]{\mathcal{W}(v_1)} E v_3 C v_1
 \end{array}$$

It works, but the effect is so limited.

Next Approach: Reordering Statements and Arguments

No side effects on SLPs; thus, we can reorder statements and arguments.

$v_1 \leftarrow A \oplus B;$	# ₁ / ₀ [5]		$v_2 \leftarrow \oplus(A, D, E);$	# ₁ / ₀ [5]
$v_2 \leftarrow \oplus(E, D, A);$	[7]	Reordering →	$v_1 \leftarrow A \oplus B;$	[3]
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$	[3]
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_4 \leftarrow v_1 \oplus C;$	[3]
$\text{return}(v_2, v_3, v_4);$	20		$\text{return}(v_2, v_3, v_4);$	14

Next Approach: Reordering Statements and Arguments

No side effects on SLPs; thus, we can reorder statements and arguments.

$v_1 \leftarrow A \oplus B;$	#1/0 [5]		$v_2 \leftarrow \oplus(A, D, E);$	#1/0 [5]
$v_2 \leftarrow \oplus(E, D, A);$	[7]	Reordering \Rightarrow	$v_1 \leftarrow A \oplus B;$	[3]
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$	[3]
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_4 \leftarrow v_1 \oplus C;$	[3]
$\text{return}(v_2, v_3, v_4);$	20		$\text{return}(v_2, v_3, v_4);$	14

Using *Pebble Game*, we can integrate $\left\{ \begin{array}{l} \text{Recycling Variables and} \\ \text{Reordering} \end{array} \right.$

★ R. Sethi, 1975, *Complete register allocation problems*.

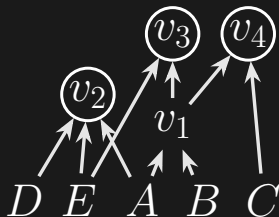
Next Approach: Reordering Statements and Arguments

No side effects on SLPs; thus, we can reorder statements and arguments.

	# ₁ / ₀		# ₁ / ₀
$v_1 \leftarrow A \oplus B;$	[5]	$\xrightarrow{\text{Reordering}}$	$v_2 \leftarrow \oplus(A, D, E);$ [5]
$v_2 \leftarrow \oplus(E, D, A);$	[7]		$v_1 \leftarrow A \oplus B;$ [3]
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$ [3]
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_4 \leftarrow v_1 \oplus C;$ [3]
$\text{return}(v_2, v_3, v_4);$	20		$\text{return}(v_2, v_3, v_4);$ 14

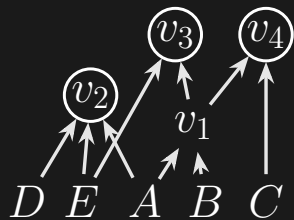
Using *Pebble Game*, we can integrate $\left\{ \begin{array}{l} \text{Recycling Variables and} \\ \text{Reordering} \end{array} \right.$

- ★ R. Sethi, 1975, *Complete register allocation problems*.
- ▶ We play the pebble game on DAGs or abstract syntax graphs.
- ▶ We aim to put pebbles in return nodes.



Pebble Game & Intractability of Optimization Problem

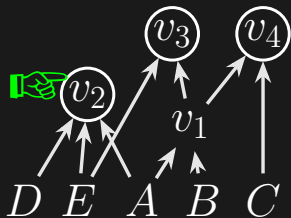
Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



Example: Evaluating strategy based on Depth-first-search

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



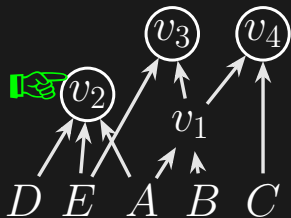
v_2 :

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



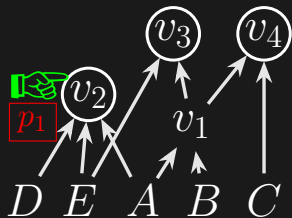
$$v_2 : \quad \leftarrow \oplus(A, D, E);$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



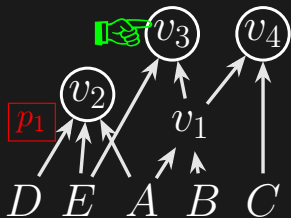
$$v_2 : p_1 \leftarrow \oplus(A, D, E);$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.
3. Put a pebble p_1 on v_2 to denote v_2 is visited.

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



$$v_2 : p_1 \leftarrow \oplus(A, D, E);$$

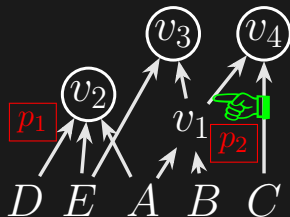
$$v_3 : \leftarrow E \oplus$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.
3. Put a pebble p_1 on v_2 to denote v_2 is visited.
4. Choose v_3 from 2 unvisited roots, and first visit E .

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



$$v_2 : p_1 \leftarrow \oplus(A, D, E);$$

$$v_1 : p_2 \leftarrow A \oplus B;$$

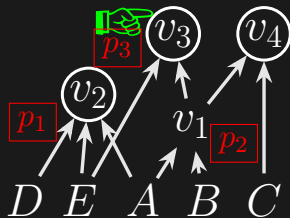
$$v_3 : \leftarrow E \oplus$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.
3. Put a pebble p_1 on v_2 to denote v_2 is visited.
4. Choose v_3 from 2 unvisited roots, and first visit E .
5. Visit the unvisited child v_1 of v_3 , evaluate, and pebble p_2

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



$$v_2 : p_1 \leftarrow \oplus(A, D, E);$$

$$v_1 : p_2 \leftarrow A \oplus B;$$

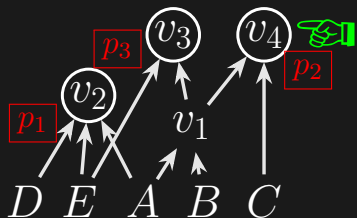
$$v_3 : p_3 \leftarrow E \oplus p_2;$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.
3. Put a pebble p_1 on v_2 to denote v_2 is visited.
4. Choose v_3 from 2 unvisited roots, and first visit E .
5. Visit the unvisited child v_1 of v_3 , evaluate, and pebble p_2
6. Back to v_3 and pebble p_3

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



$$v_2 : p_1 \leftarrow \oplus(A, D, E);$$

$$v_1 : p_2 \leftarrow A \oplus B;$$

$$v_3 : p_3 \leftarrow E \oplus p_2;$$

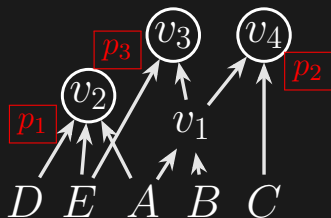
$$v_4 : p_2 \leftarrow C \oplus p_2;$$

Example: Evaluating strategy based on Depth-first-search

1. Choose v_2 from unvisited roots: alphabetical small $v_2 \prec v_3 \prec v_4$.
2. Evaluate the children of v_2 in alphabetical order.
3. Put a pebble p_1 on v_2 to denote v_2 is visited.
4. Choose v_3 from 2 unvisited roots, and first visit E .
5. Visit the unvisited child v_1 of v_3 , evaluate, and pebble p_2
6. Back to v_3 and pebble p_3
7. Finally, we compute v_4 with *moving/recycling* pebble p_2 .

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling $\#_{I/O}$



```
v2 : p1 ← ⊕(A, D, E); [7]
v1 : p2 ← A ⊕ B; [3]
v3 : p3 ← E ⊕ p2; [3]
v4 : p2 ← C ⊕ p2; [2]
return(p1, p3, p2); 15
```

Example: Evaluating strategy based on Depth-first-search

Can we find the best reordering and pebbling in polynomial time?

Theorem (Sethi 1975, Papp & Wattenhofer 2020)

Unless $\mathbf{P} = \mathbf{NP}$, for a given P , in polynomial time, we cannot find a Q that $\llbracket P \rrbracket = \llbracket Q \rrbracket$ and minimizes $\#_{I/O}(Q)$.

We use DFS-based strategy as above in our evaluation.

Evaluation

Data Set & Evaluation Environment

We consider RS(10, 4) as an example data set.

- ▶ We have 1-encoding SLP P_{enc} .
- ▶ We have $\binom{14}{4} = 1001$ decoding SLPs.

We used two environments in my paper:

name	CPU	Clock	Core	RAM
intel	i7-7567U	4.0GHz	2	DDR3-2133 16GB
amd	Ryzen 2600	3.9GHz	6	DDR4-2666 48GB

In a distributed computation,
our test environments correspond to single nodes.

L1 cache specification:	Size	Associativity	Line Size
	32KB/core	8-way	64 bytes

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755			
$\#_{mem}$	2265			

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755			
$\#_{\text{mem}}$	2265			
$\mathcal{B} = 512 : \#_{\text{I/O}}(K = 64)$	570			
$\mathcal{B} = 1\text{K} : \#_{\text{I/O}}(K = 32)$	1262			
$\mathcal{B} = 2\text{K} : \#_{\text{I/O}}(K = 16)$	1598			

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

\mathcal{B} -Byte Blocking for Cache Efficiency

```
for  $i \leftarrow 0 \dots (A.\text{len} / \mathcal{B})$  {  
   $v_1 = \text{xor}(A, B);$             $v_1^{[i]} = \text{xor}(A^{[i]}, B^{[i]});$   
   $v_2 = \text{xor}(v_1, C, D);$   $\implies$   $v_2^{[i]} = \text{xor}(v_1^{[i]}, C^{[i]}, D^{[i]});$   
   $\text{return}(v_1, v_2);$        }  
   $\text{return}(v_1, v_2);$ 
```

where $A^{[i]}$ is the i -th \mathcal{B} -bytes block.

$\mathcal{B} = 2K : \#_{\text{I/O}}(K = 16)$

1598

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755			
$\#_{mem}$	2265			
$\mathcal{B} = 512 :$				
$\#_{I/O}(K = 64)$	570			
Throughput (GB/s)	3.10			
$\mathcal{B} = 1K :$				
$\#_{I/O}(K = 32)$	1262			
Throughput (GB/s)	4.03			
$\mathcal{B} = 2K :$				
$\#_{I/O}(K = 16)$	1598			
Throughput (GB/s)	4.45			

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Dabbling
$\#_{\oplus}$	755			
$\#_{mem}$	2265			
$\mathcal{B} = 512 :$ $\#_{I/O}(K = 64)$	570			
Throughput (GB/s)	3.10			
$\mathcal{B} = 1K :$ $\#_{I/O}(K = 32)$	1262			
Throughput (GB/s)	4.03			
$\mathcal{B} = 2K :$ $\#_{I/O}(K = 16)$	1598			
Throughput (GB/s)	4.45			

Why smaller blocks are slower than the large one?

Pros: Smaller blocks,

- ▶ More cache-able blocks $\frac{32K}{\mathcal{B}}$.

Cons: Smaller blocks,

- ▶ Due to cache conflicts, using cache identically is more difficult.
- ▶ Latency penalty becomes totally large.

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#\oplus$	755	385		
$\#\text{mem}$	2265	1155		
$\mathcal{B} = 512 :$				
$\#\text{I/O}(K = 64)$	570			
Throughput (GB/s)	3.10			
$\mathcal{B} = 1\text{K} :$				
$\#\text{I/O}(K = 32)$	1262			
Throughput (GB/s)	4.03			
$\mathcal{B} = 2\text{K} :$				
$\#\text{I/O}(K = 16)$	1598			
Throughput (GB/s)	4.45			

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755	385		
$\#_{\text{mem}}$	2265	1155		
$\mathcal{B} = 512 :$ $\#_{\text{I/O}}(K = 64)$	570	1231		
Throughput (GB/s)	3.10			
$\mathcal{B} = 1\text{K} :$ $\#_{\text{I/O}}(K = 32)$	1262	1465		
Throughput (GB/s)	4.03			
$\mathcal{B} = 2\text{K} :$ $\#_{\text{I/O}}(K = 16)$	1598	1599		
Throughput (GB/s)	4.45			

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755	385		
$\#_{\text{mem}}$	2265	1155		
$\mathcal{B} = 512 :$				
$\#_{\text{I/O}}(K = 64)$	570	1231		
Throughput (GB/s)	3.10	4.18		
$\mathcal{B} = 1\text{K} :$				
$\#_{\text{I/O}}(K = 32)$	1262	1465		
Throughput (GB/s)	4.03	4.36		
$\mathcal{B} = 2\text{K} :$				
$\#_{\text{I/O}}(K = 16)$	1598	1599		
Throughput (GB/s)	4.45	4.86		

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#\oplus$	755	385	N/A	
$\#\text{mem}$	2265	1155	677	
$\mathcal{B} = 512 :$				
$\#\text{I/O}(K = 64)$	570	1231		
Throughput (GB/s)	3.10	4.18		
$\mathcal{B} = 1\text{K} :$				
$\#\text{I/O}(K = 32)$	1262	1465		
Throughput (GB/s)	4.03	4.36		
$\mathcal{B} = 2\text{K} :$				
$\#\text{I/O}(K = 16)$	1598	1599		
Throughput (GB/s)	4.45	4.86		

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#_{\oplus}$	755	385	N/A	
$\#_{\text{mem}}$	2265	1155	677	
$\mathcal{B} = 512 :$				
$\#_{\text{I/O}}(K = 64)$	570	1231	936	
Throughput (GB/s)	3.10	4.18		
$\mathcal{B} = 1\text{K} :$				
$\#_{\text{I/O}}(K = 32)$	1262	1465	1086	
Throughput (GB/s)	4.03	4.36		
$\mathcal{B} = 2\text{K} :$				
$\#_{\text{I/O}}(K = 16)$	1598	1599	1144	
Throughput (GB/s)	4.45	4.86		

Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#\oplus$	755	385	N/A	
$\#\text{mem}$	2265	1155	677	
$\mathcal{B} = 512 :$				
$\#\text{I/O}(K = 64)$	570	1231	936	
Throughput (GB/s)	3.10	4.18	6.98	
$\mathcal{B} = 1\text{K} :$				
$\#\text{I/O}(K = 32)$	1262	1465	1086	
Throughput (GB/s)	4.03	4.36	7.50	
$\mathcal{B} = 2\text{K} :$				
$\#\text{I/O}(K = 16)$	1598	1599	1144	
Throughput (GB/s)	4.45	4.86	7.12	

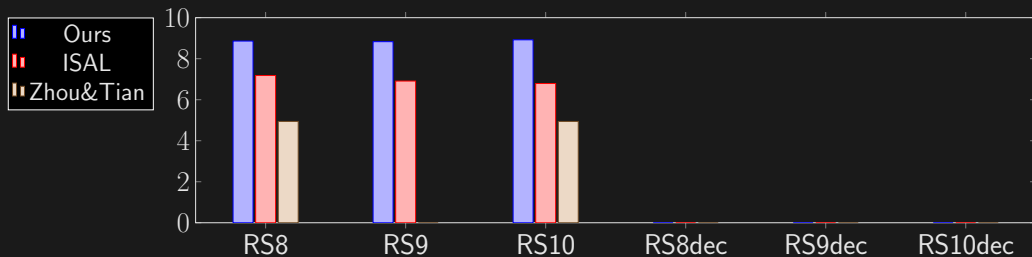
Improvements by heuristics for the encoding SLP on Intel PC

Throughput is Avg. of 1000-runs for 10MB randomly generated data

Metric	Base P_{enc}	RePair	RePair + Fuse	RePair + Fuse + Pebbling
$\#\oplus$	755	385	N/A	
$\#\text{mem}$	2265	1155	677	
$\mathcal{B} = 512 :$				
$\#\text{I/O}(K = 64)$	570	1231	936	636
Throughput (GB/s)	3.10	4.18	6.98	7.24
$\mathcal{B} = 1\text{K} :$				
$\#\text{I/O}(K = 32)$	1262	1465	1086	779
Throughput (GB/s)	4.03	4.36	7.50	8.92
$\mathcal{B} = 2\text{K} :$				
$\#\text{I/O}(K = 16)$	1598	1599	1144	845
Throughput (GB/s)	4.45	4.86	7.12	8.55

Throughput Comparison (Intel + 1K-Blocking)

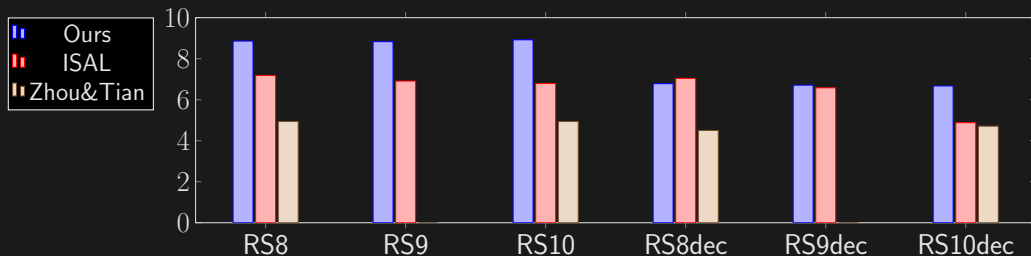
Enc	#mem	#I/O	Ours	ISA-L v2.30	Zhou & Tian
RS(8, 4)	543	585	8.86 GB/s	7.18 GB/s	4.94 GB/s
RS(9, 4)	611	671	8.83	6.91	N/A in their paper
RS(10, 4)	677	779	8.92	6.79	4.94



Throughput Comparison (Intel + 1K-Blocking)

Enc	#mem	#I/O	Ours	ISA-L v2.30	Zhou & Tian
RS(8, 4)	543	585	8.86 GB/s	7.18 GB/s	4.94 GB/s
RS(9, 4)	611	671	8.83	6.91	N/A in their paper
RS(10, 4)	677	779	8.92	6.79	4.94

Dec	#mem	#I/O	Ours	ISA-L v2.30	Zhou & Tian
RS(8, 4)	747	811	6.78 GB/s	7.04 GB/s	4.50 GB/s
RS(9, 4)	829	968	6.71	6.58	N/A
RS(10, 4)	923	1077	6.67	4.88	4.71



Conclusion (+ Other Throughput Scores)

intel 1K (GB/sec)	Ours		ISA-L v 2.30		Zhou & Tian	
	Enc	Dec	Enc	Dec	Enc	Dec
RS (8, 3)	12.32	8.82	9.09	9.25	6.08	5.57
RS (9, 3)	11.97	8.27	7.31	7.92	6.17	5.66
RS (10, 3)	11.78	8.89	6.78	7.93	6.15 _S	5.90
RS (8, 2)	18.79	14.59	12.99	13.34	8.13 _E	8.07 _E
RS (9, 2)	18.93	14.27	11.85	12.03	8.34 _E	8.04
RS (10, 2)	18.98	14.66	12.12	12.61	8.40 _E	8.22 _E

Conclusion

- ▶ We identified bitmatrix multiplication as straight line programs (SLP).
- ▶ We optimized XOR-based EC by optimizing SLPs using various program optimization techniques.
- ▶ Each of our techniques is not difficult; however, it suffices to match Intel's high performance library ISAL.
- ▶ As future work on cache optimization, I plan to accommodate multi-layer cache L1, L2, and L3 cache.