

プログラムの最適化手法を用いた Erasure Coding の最適化

上里 友弥 @ Cyber Agent

2024 / 02 / 11 @ 筑波大学

今日のお話: 非正方行列積を、爆速で実行したい

$$14 \underbrace{\begin{bmatrix} x_{1,1} & \cdots & x_{1,10} \\ x_{2,1} & \cdots & x_{2,10} \\ \vdots & \ddots & \vdots \\ x_{12,1} & \cdots & x_{12,10} \\ x_{13,1} & \cdots & x_{13,10} \\ x_{14,1} & \cdots & x_{14,10} \end{bmatrix}}_{10} \times \underbrace{\begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,1000000} & \cdots & y_{1,4\text{百万}} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,1000000} & \cdots & y_{2,4\text{百万}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{10,1} & y_{10,2} & \cdots & y_{10,1000000} & \cdots & y_{10,4\text{百万}} \end{bmatrix}}_{\text{かなりナゲエ行列}}$$

今日のお話: 非正方行列積を、爆速で実行したい

$$14 \underbrace{\begin{bmatrix} x_{1,1} & \cdots & x_{1,10} \\ x_{2,1} & \cdots & x_{2,10} \\ \vdots & \ddots & \vdots \\ x_{12,1} & \cdots & x_{12,10} \\ x_{13,1} & \cdots & x_{13,10} \\ x_{14,1} & \cdots & x_{14,10} \end{bmatrix}}_{10} \times \underbrace{\begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,1000000} & \cdots & y_{1,4\text{百万}} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,1000000} & \cdots & y_{2,4\text{百万}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{10,1} & y_{10,2} & \cdots & y_{10,1000000} & \cdots & y_{10,4\text{百万}} \end{bmatrix}}_{\text{かなりナゲエ行列}}$$

→ 速度がユーザー体験に直結する状況から来た問題; コマケエこた後述

今日のお話: 非正方行列積を、爆速で実行したい

$$14 \underbrace{\begin{bmatrix} x_{1,1} & \cdots & x_{1,10} \\ x_{2,1} & \cdots & x_{2,10} \\ \vdots & \ddots & \vdots \\ x_{12,1} & \cdots & x_{12,10} \\ x_{13,1} & \cdots & x_{13,10} \\ x_{14,1} & \cdots & x_{14,10} \end{bmatrix}}_{10} \times \underbrace{\begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,1000000} & \cdots & y_{1,4\text{百万}} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,1000000} & \cdots & y_{2,4\text{百万}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{10,1} & y_{10,2} & \cdots & y_{10,1000000} & \cdots & y_{10,4\text{百万}} \end{bmatrix}}_{\text{かなりナゲエ行列}}$$

- 速度がユーザー体験に直結する状況から来た問題; コマケエコた後述
- ただし x や y は $\mathbb{F}[2^8]$ という 有限体 の要素
有理数体 \mathbb{Q} , 実数体 \mathbb{R} , 複素数体 \mathbb{C} ではない
- $\mathbb{F}[2^8]$ の積の実行 $\mathbb{F}[2^8] \ni x, y \mapsto x \cdot y$ は「マジ遅い」

今日のお話: 非正方行列積を、爆速で実行したい

$$14 \underbrace{\begin{bmatrix} x_{1,1} & \cdots & x_{1,10} \\ x_{2,1} & \cdots & x_{2,10} \\ \vdots & \ddots & \vdots \\ x_{12,1} & \cdots & x_{12,10} \\ x_{13,1} & \cdots & x_{13,10} \\ x_{14,1} & \cdots & x_{14,10} \end{bmatrix}}_{10} \times \underbrace{\begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,1000000} & \cdots & y_{1,4\text{百万}} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,1000000} & \cdots & y_{2,4\text{百万}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{10,1} & y_{10,2} & \cdots & y_{10,1000000} & \cdots & y_{10,4\text{百万}} \end{bmatrix}}_{\text{かなりナゲエ行列}}$$

- 速度がユーザー体験に直結する状況から来た問題; コマケエこた後述
- ただし x や y は $\mathbb{F}[2^8]$ という 有限体 の要素
有理数体 \mathbb{Q} , 実数体 \mathbb{R} , 複素数体 \mathbb{C} ではない
- $\mathbb{F}[2^8]$ の積の実行 $\mathbb{F}[2^8] \ni x, y \mapsto x \cdot y$ は「マジ遅い」
- 「マジ遅い」演算を $(10 \cdot 4\text{百万}) \cdot 14$ 回実行してたら終わる

今日のお話: 非正方行列積を、爆速で実行したい

$$14 \underbrace{\begin{bmatrix} x_{1,1} & \cdots & x_{1,10} \\ x_{2,1} & \cdots & x_{2,10} \\ \vdots & \ddots & \vdots \\ x_{12,1} & \cdots & x_{12,10} \\ x_{13,1} & \cdots & x_{13,10} \\ x_{14,1} & \cdots & x_{14,10} \end{bmatrix}}_{10} \times \underbrace{\begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,1000000} & \cdots & y_{1,4\text{百万}} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,1000000} & \cdots & y_{2,4\text{百万}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{10,1} & y_{10,2} & \cdots & y_{10,1000000} & \cdots & y_{10,4\text{百万}} \end{bmatrix}}_{\text{かなりナゲエ行列}}$$

- 速度がユーザー体験に直結する状況から来た問題; コマケエこた後述
- ただし x や y は $\mathbb{F}[2^8]$ という 有限体 の要素
有理数体 \mathbb{Q} , 実数体 \mathbb{R} , 複素数体 \mathbb{C} ではない
- $\mathbb{F}[2^8]$ の積の実行 $\mathbb{F}[2^8] \ni x, y \mapsto x \cdot y$ は「マジ遅い」
- 「マジ遅い」演算を $(10 \cdot 4\text{百万}) \cdot 14$ 回実行してたら終わる

なんとかする。プログラム理論のプロとして——
(有限体のプロでも、行列のプロでもないんで 😊)



Accelerating XOR-Based Erasure Coding using Program Optimization Techniques

Yuya Uezato
Dwango, Co., Ltd.
Japan
yuuya_uezato@dwango.co.jp

ABSTRACT

Erasure coding (EC) affords data redundancy for large-scale systems. XOR-based EC is an easy-to-implement method for optimizing EC. This paper addresses a significant performance gap between the state-of-the-art XOR-based EC approach (~4.9 GB/s coding throughput) and Intel's high-performance EC library based on another approach (~6.7 GB/s). We propose a novel approach based on our observation that XOR-based EC virtually generates programs of a Domain Specific Language for XORing byte arrays. We formalize such programs as straight-line programs (SLPs) of compiler construction and optimize SLPs using various program optimization techniques. Our optimization flow is three-fold: 1) reducing the number of XORs using grammar compression algorithms; 2) reducing memory accesses using deforestation, a functional program optimization method; and 3) reducing cache misses using the (red-blue) pebble game of program analysis. We provide an experimental library, which outperforms Intel's library with an ~8.92 GB/s throughput.

CCS CONCEPTS

- Software and its engineering → Compilers; Context specific languages;
- Mathematics of computing → Coding theory;
- Computer systems organization → Embedded systems

can store 10-times more objects than through replication; however, we cannot recover data if five nodes are down. Another distributed system Ceph [24] offers $\text{RS}(n, p)$ for any n and p . On Linux, we can use RAID-6, a codec similar to $\text{RS}(n, 2)$ [11, 82]. Using EC instead of replication degrades the system performance since the encoding and decoding of EC are heavy computation and are required for each storing to and loading from a system. It is often stated that EC is suitable only for archiving cold (rarely accessed) data [29, 49, 93].

We clarify the pros and cons of EC by observing how RS works. To encode data using matrix multiplication (hereafter we use the acronym MM), RS adopts matrices over \mathbb{F}_{2^8} , the finite field with $2^8 = 256$ elements. Since each element of \mathbb{F}_{2^8} is coded by one byte (8 bits), we can identify an N -bytes data as an N -elements array of \mathbb{F}_{2^8} . $\text{RS}(n, p)$ encodes an N -bytes data D using an $(n + p) \times n$ Vandermonde matrix $\mathcal{V} \in \mathbb{F}_{2^8}^{(n+p) \times n}$, which is crucial for decoding as we will see below, as follows:

$$\begin{matrix} & \begin{pmatrix} \vec{d}_1 \\ \vdots \\ \vec{d}_n \\ \vdots \\ \vec{d}_{n+p} \end{pmatrix} \\ \begin{pmatrix} n \\ + \\ p \end{pmatrix} \times \mathcal{V} \in \mathbb{F}_{2^8}^{(n+p) \times n} & = \begin{pmatrix} \vec{b}_1 \\ \vdots \\ \vec{b}_n \\ \vdots \\ \vec{b}_{n+p} \end{pmatrix} \end{matrix} \quad \text{where}$$

\mathcal{V} is the MM over \mathbb{F}_{2^8} ;
 \vec{d}_i is i -th $\frac{N}{n}$ -bytes block of D ;
 \vec{b}_j is an $\frac{N}{n}$ -bytes coded block.



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:

$$\begin{array}{l} \text{(for Standard EC)} \quad A \times B \text{ over } \mathbb{F}[2^8], \\ \text{(for XOR-based EC)} \quad C \times D \text{ over } \mathbb{F}[2]. \end{array}$$

- ▶ *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 ...

Optimizing

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

$\mathbb{F}[2^8]$ 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.
- ▶ *Byte Field* (We will see it in the later page).
- ▶ *Bit Field*

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

Optimizing

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

$\mathbb{F}[2^8]$ 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.
- ▶ *Byte Level* (We will see it in the later page).
- ▶ *Bit Field*

$\mathbb{F}[2] = \{0, 1\}$ 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:

$$\begin{array}{ll} \text{(for Standard EC)} & A \times B \text{ over } \mathbb{F}[2^8], \\ \text{(for XOR-based EC)} & C \times D \text{ over } \mathbb{F}[2]. \end{array}$$

- ▶ *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:

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

(for XOR-based EC) $C \times D$ 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:

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

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

Optimizing matrix multiplication over two finite fields:

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

(for XOR-based EC) $C \times D$ 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.

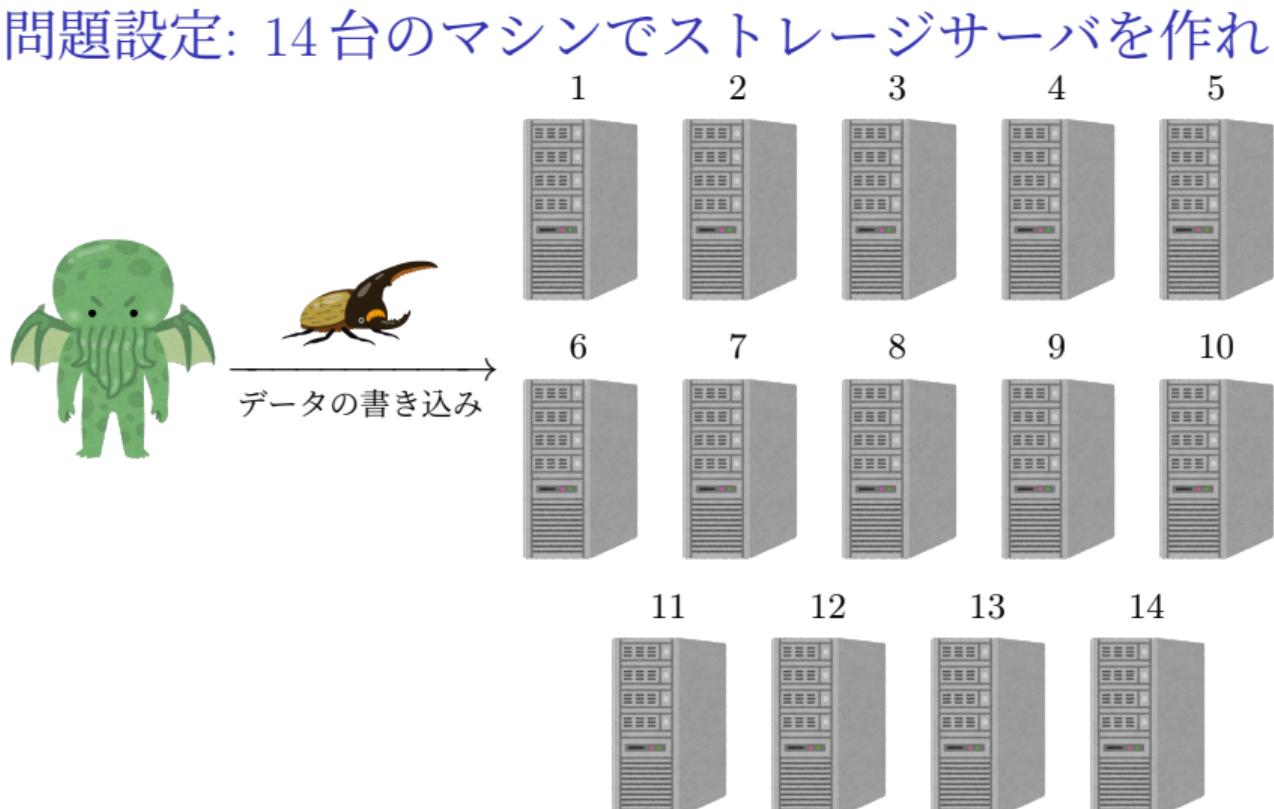
A is small. B is large. In this talk, as an example, we consider:

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

This setting comes from *erasure coding*.

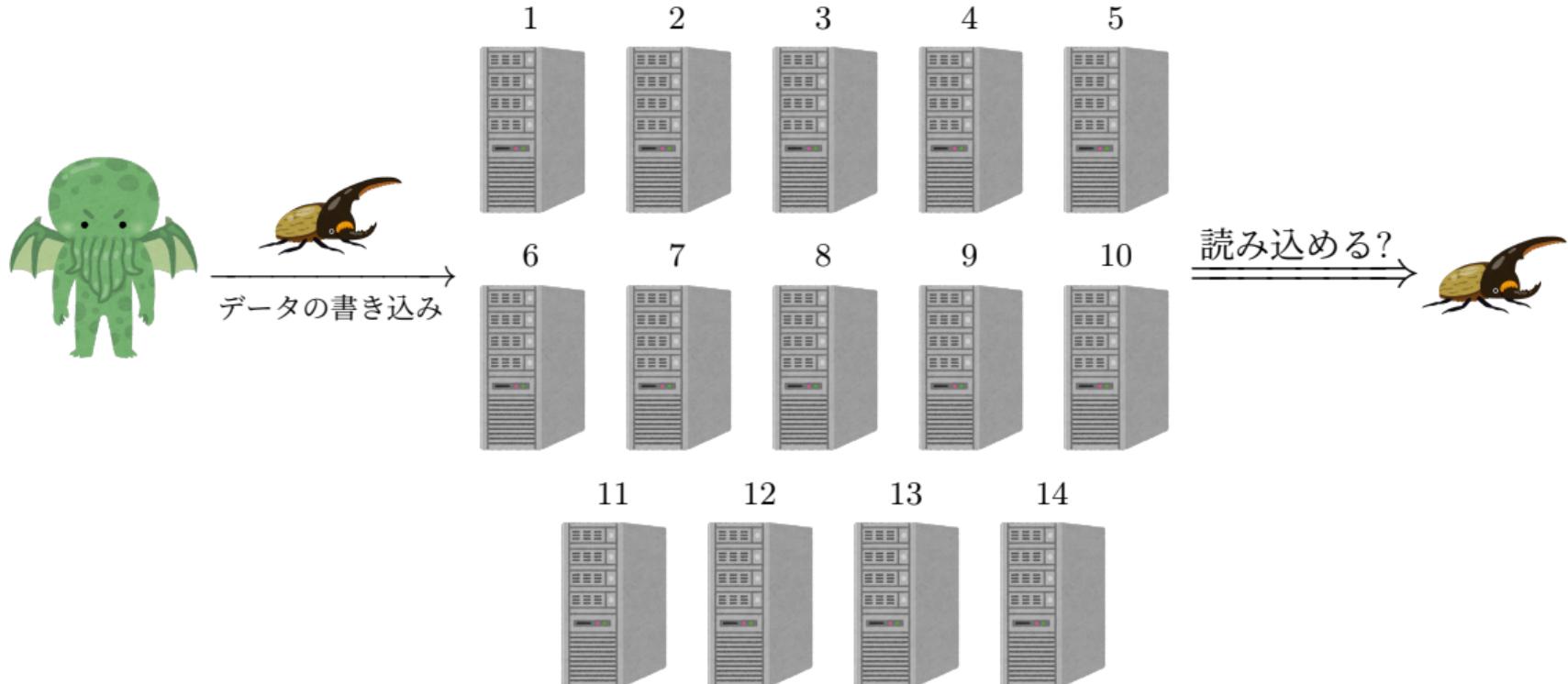
What are Erasure Coding (EC) and XOR-Based EC

問題設定: 14台のマシンでストレージサーバを作れ



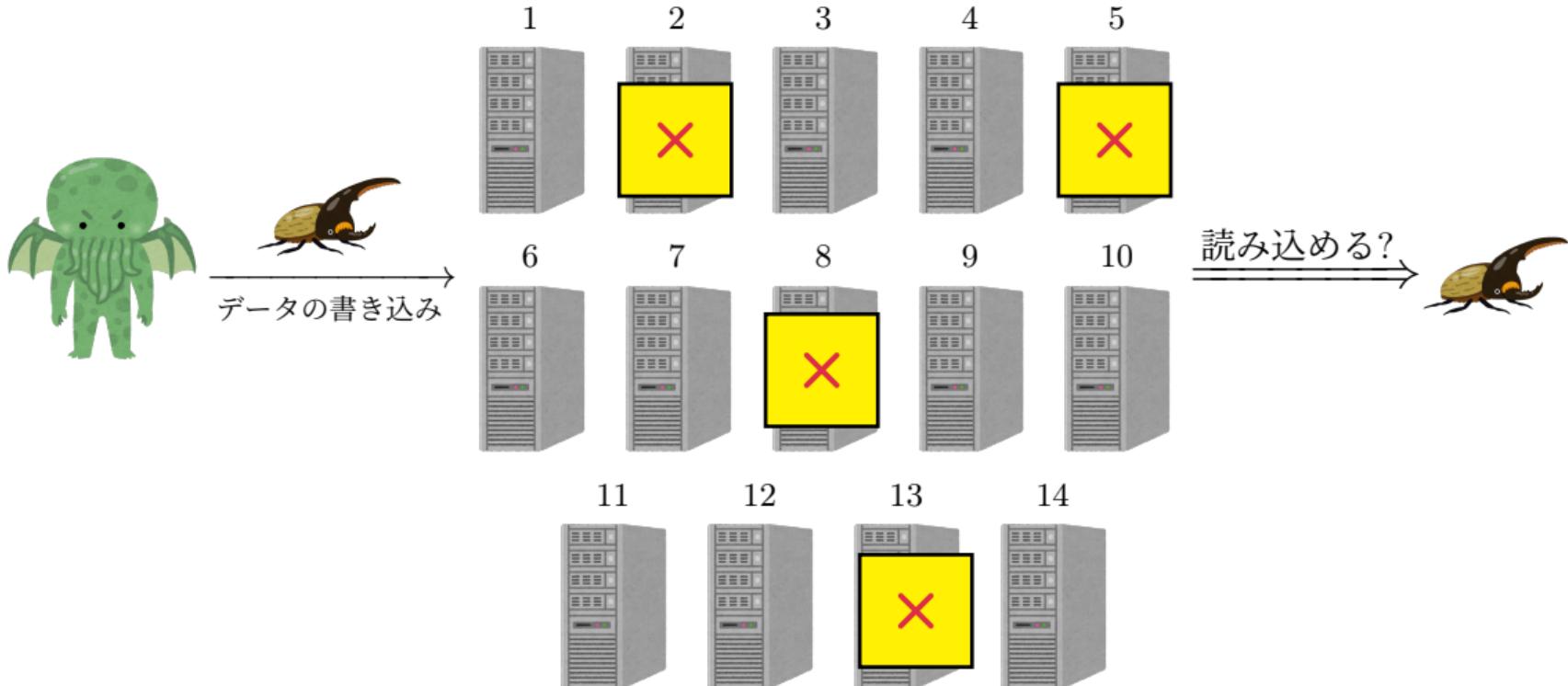
※ 一台の容量は 20TB;

問題設定: 14台のマシンでストレージサーバを作れ



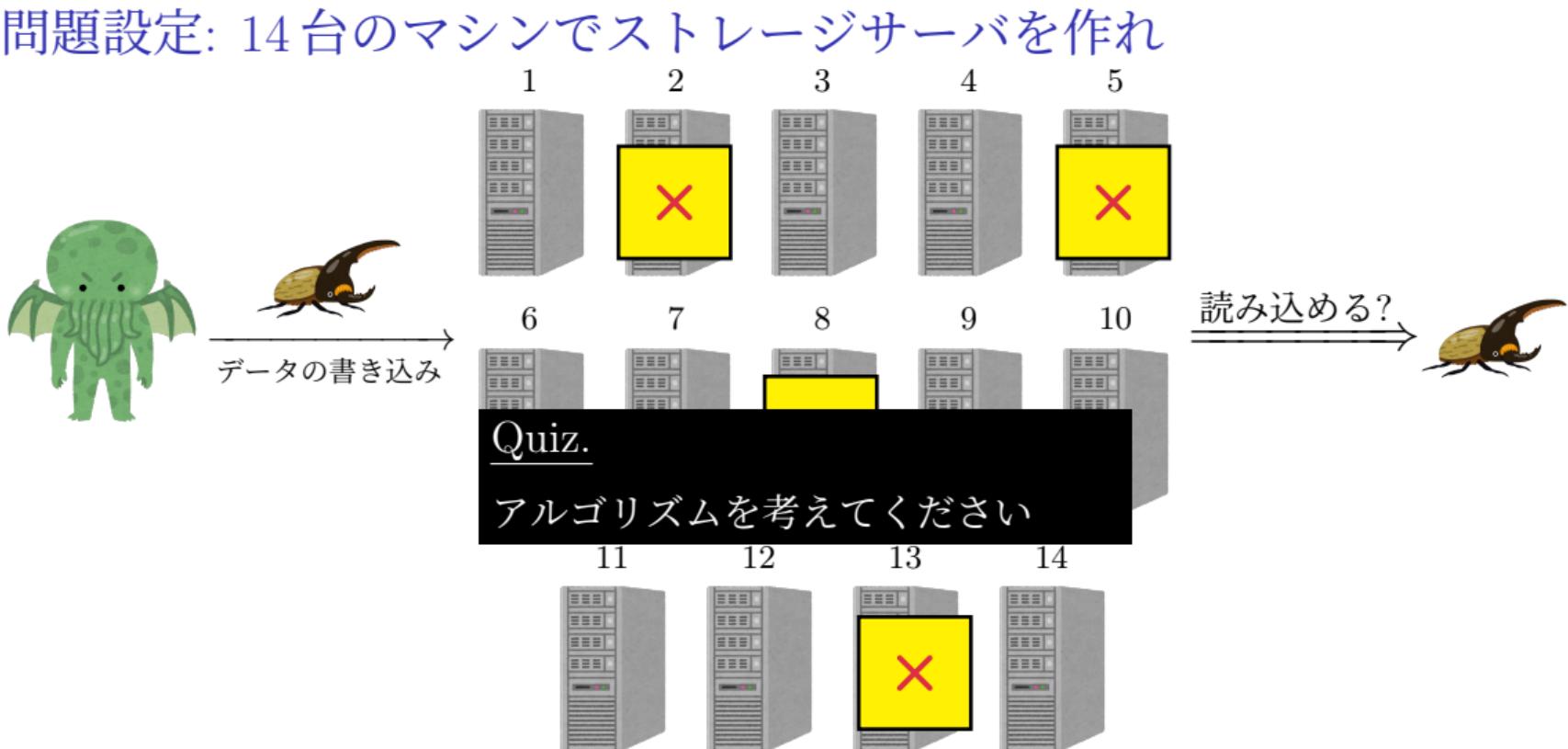
※ 一台の容量は 20TB;

問題設定: 14台のマシンでストレージサーバを作れ



- 一台の容量は 20TB; システム容量は **200TB** とせよ
- 四台壊れてても データ復元しる! (五台壊れてたら諦める)

問題設定: 14台のマシンでストレージサーバを作れ



- 一台の容量は 20TB; システム容量は **200TB** とせよ
- 四台壊れてても データ復元し！ (五台壊れてたら諦める)

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.

試しに RAID を使ってみる

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.

試しに RAID を使ってみる

RAID-1 別名 Replication: ユーザからの入力を全台にコピーする。

- ▶ 冗長性最強: 13台壊れても大丈夫
- ▶ 空間効率悪し: 結局一台分なので今回は 20TB しかない

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.

試しに RAID を使ってみる

RAID-1 別名 Replication: ユーザからの入力を全台にコピーする。

- ▶ 冗長性最強: 13台壊れても大丈夫
- ▶ 空間効率悪し: 結局一台分なので今回は 20TB しかない

RAID-5 データを 13分割し、parity と呼ばれるものを 1つ作る。

- ▶ 冗長性は一応ある: 1台だけなら壊れても大丈夫
- ▶ 空間効率: $13 \times 20\text{TB} = 260\text{TB}$ の総容量

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.

試しに RAID を使ってみる

RAID-1 別名 Replication: ユーザからの入力を全台にコピーする。

- ▶ 冗長性最強: 13台壊れても大丈夫
- ▶ 空間効率悪し: 結局一台分なので今回は 20TB しかない

RAID-5 データを 13分割し、parity と呼ばれるものを 1つ作る。

- ▶ 冗長性は一応ある: 1台だけなら壊れても大丈夫
- ▶ 空間効率: $13 \times 20\text{TB} = 260\text{TB}$ の総容量

RAID-6 データを 12分割し、parity を 2つ作る。

- ▶ 冗長性はやや良い: 2台なら壊れても大丈夫
- ▶ 空間効率: $12 \times 20\text{TB} = 240\text{TB}$ の総容量

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.
 - ▶ $14 \cdot 20 - 200 = 80$ TB can be used for data redundancy.

For this criteria, we can employ Reed-Solomon EC **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. ($p = 4$).

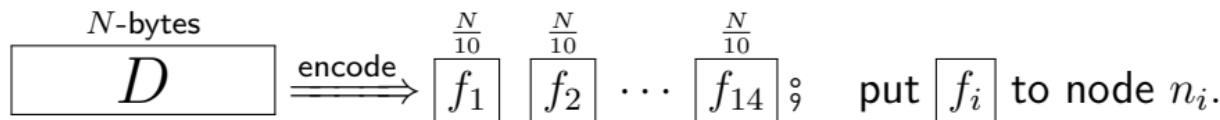
Erasure Coding: Method for Data Redundancy

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.
 - ▶ $14 \cdot 20 - 200 = 80$ TB can be used for data redundancy.

For this criteria, we can employ Reed-Solomon EC **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. ($p = 4$).



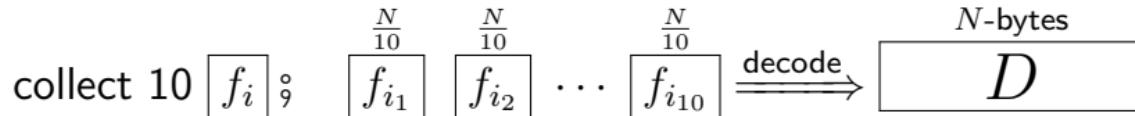
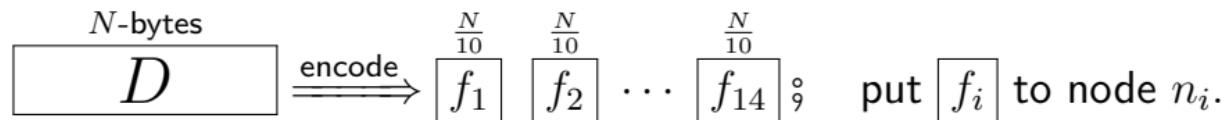
Erasure Coding: Method for Data Redundancy

Example (Building a streaming media server with criteria)

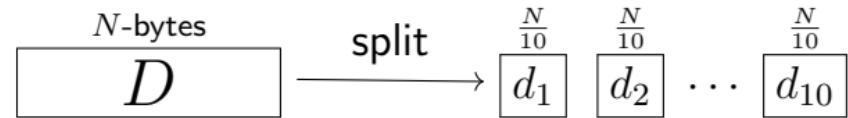
1. We have 14 nodes. Each node has a 20TB disk.
2. We can load data even if nodes ≤ 4 are down.
3. The total capacity of our server = 200TB.
 - ▶ $14 \cdot 20 - 200 = 80$ TB can be used for data redundancy.

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

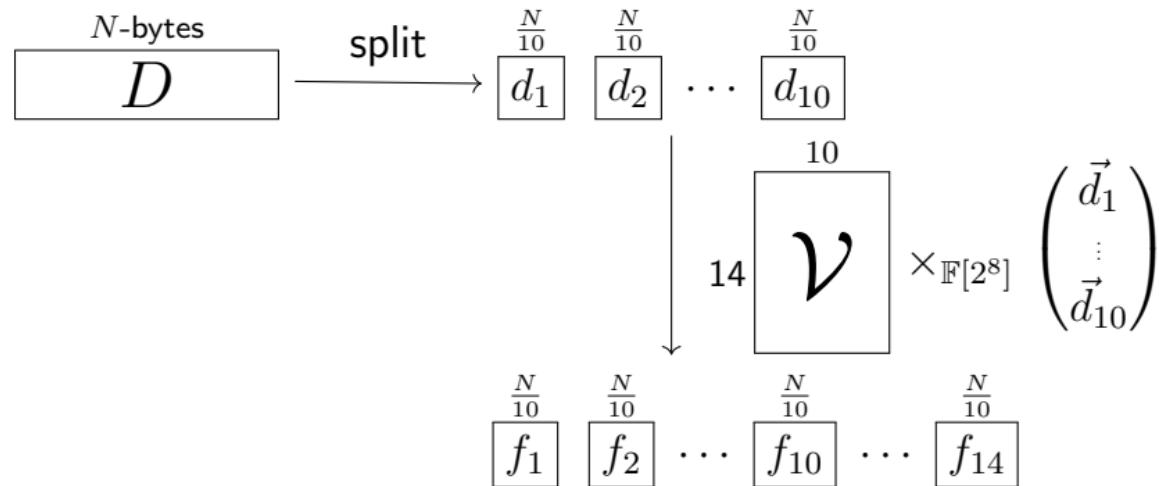
- ▶ d : we can assume d -nodes are living. ($d = 14 - p = 10$).
- ▶ p : we can permit nodes $\leq p$ go down. ($p = 4$).



How encoding and decoding are implemented in RS(10, 4) ?

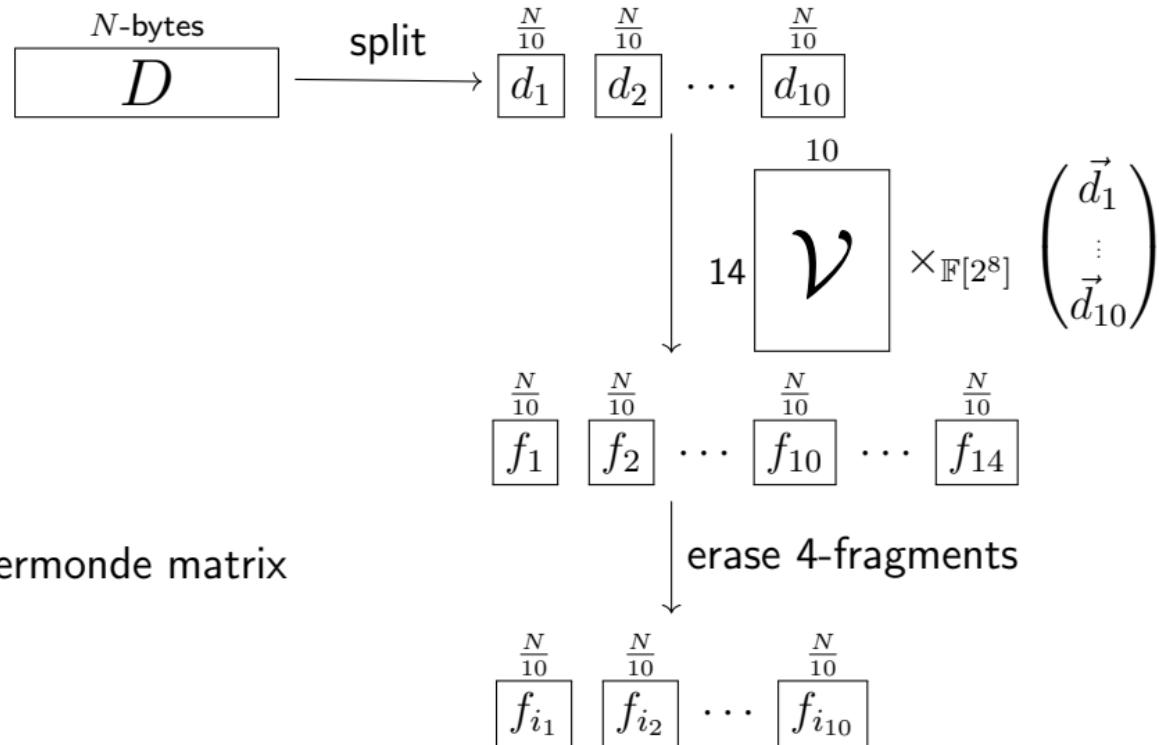


How encoding and decoding are implemented in RS(10, 4) ?

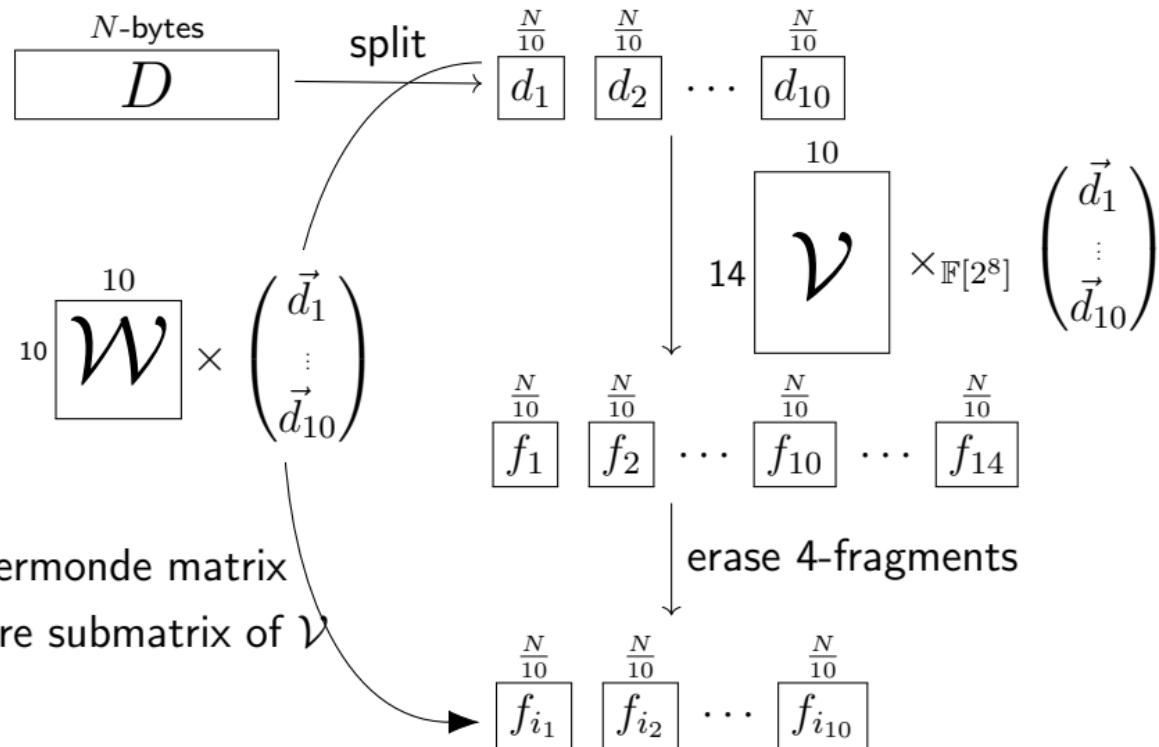


- ▶ \mathcal{V} : Vandermonde matrix

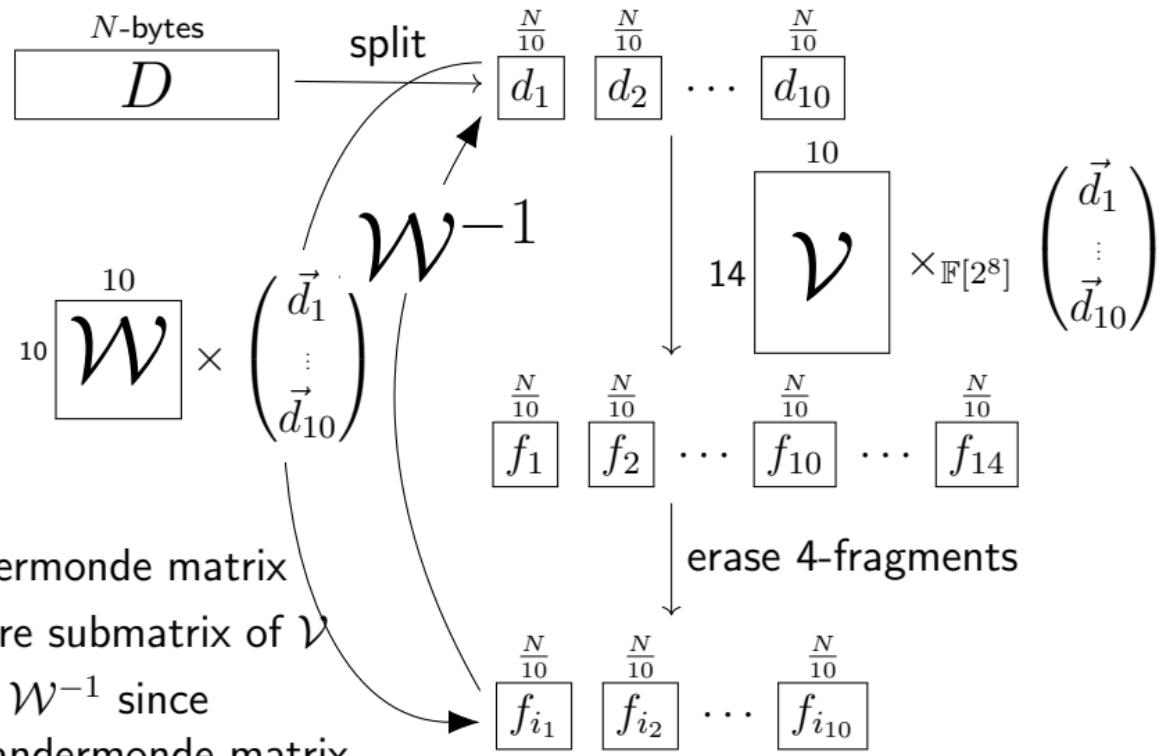
How encoding and decoding are implemented in RS(10, 4) ?



How encoding and decoding are implemented in RS(10, 4) ?



How encoding and decoding are implemented in RS(10, 4) ?



小さい行列で計算してみる

生成行列

$$\begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \\ \vec{x}_4 \\ \vec{x}_5 \end{pmatrix} \times \begin{pmatrix} & & & \\ \text{データ行列 } \mathcal{D} & & & \\ & & & \end{pmatrix} =$$

$$\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}$$

data lost

$$\boxed{\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}}^{\mathcal{M}}$$

サーバ 2, 5 が応答せず、データが欠けて、実際に集まったのが \mathcal{M} という状況

小さい行列で計算してみる

生成行列

$$\begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \\ \vec{x}_4 \\ \vec{x}_5 \end{pmatrix}$$

$$\times \overbrace{\begin{pmatrix} & & & \\ & & & \\ & & & \\ \text{データ行列 } \mathcal{D} & & & \\ & & & \end{pmatrix}}^{\text{データ行列 } \mathcal{D}} =$$

$$\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}$$

data lost

$$\boxed{\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}}$$

サーバ 2, 5 が応答せず、データが欠けて、実際に集まったのが \mathcal{M} という状況

生成行列から 2 行目と 5 行目を削除した行列 \mathcal{G} をかけても同じ:

$$\overbrace{\begin{pmatrix} \vec{x}_1 \\ \vec{x}_3 \\ \vec{x}_4 \end{pmatrix}}^{\mathcal{G}}$$

$$\times \begin{pmatrix} y_1^\downarrow & y_2^\downarrow & \cdots & y_n^\downarrow \end{pmatrix} =$$

$$\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \end{pmatrix} = \mathcal{M}$$

小さい行列で計算してみる

生成行列

$$\begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \\ \vec{x}_4 \\ \vec{x}_5 \end{pmatrix} \times \overbrace{\begin{pmatrix} y_1^\downarrow & y_2^\downarrow & \cdots & y_n^\downarrow \end{pmatrix}}^{\text{データ行列 } \mathcal{D}} =$$

$$\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}$$

data lost

$$\boxed{\begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_2 \odot y_1^\downarrow & \cdots & \vec{x}_2 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \\ \vec{x}_5 \odot y_1^\downarrow & \cdots & \vec{x}_5 \odot y_n^\downarrow \end{pmatrix}}^{\mathcal{M}}$$

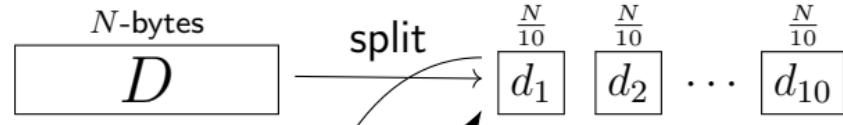
サーバ 2, 5 が応答せず、データが欠けて、実際に集まったのが \mathcal{M} という状況

生成行列から 2 行目と 5 行目を削除した行列 \mathcal{G} をかけても同じ:

$$\overbrace{\begin{pmatrix} \vec{x}_1 \\ \vec{x}_3 \\ \vec{x}_4 \end{pmatrix}}^{\mathcal{G}} \times \begin{pmatrix} y_1^\downarrow & y_2^\downarrow & \cdots & y_n^\downarrow \end{pmatrix} = \begin{pmatrix} \vec{x}_1 \odot y_1^\downarrow & \cdots & \vec{x}_1 \odot y_n^\downarrow \\ \vec{x}_3 \odot y_1^\downarrow & \cdots & \vec{x}_3 \odot y_n^\downarrow \\ \vec{x}_4 \odot y_1^\downarrow & \cdots & \vec{x}_4 \odot y_n^\downarrow \end{pmatrix} = \mathcal{M}$$

よって、逆行列があれば、 $\mathcal{G}^{-1} \times \mathcal{M} = \mathcal{G}^{-1} \times (\mathcal{G} \times \mathcal{D}) = \mathcal{D}$ で元データ \mathcal{D} が復元

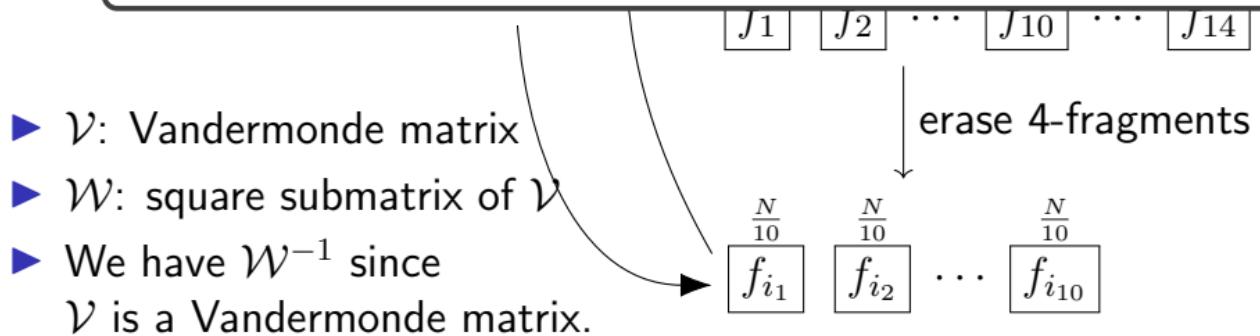
How encoding and decoding are implemented in RS(10, 4) ?



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.

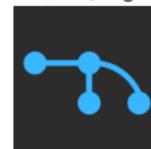


自己紹介

- ▶ 筑波大学の SCORE 研出身の博士です。博論はオートマトンのお話。

前職がドワンゴです dwango

いう 分散オブジェクトストレージを
作っていて・改良していて・基礎研究しています。



<https://github.com/frugalos/frugalos>

Frugalos

Frugal Object Storage

[crates.io](#) v1.2.0 [docs](#) failing [build](#) passing [license](#) MIT

Frugalos is a distributed object storage written by Rust.

It is suitable for storing medium size BLOBs that become petabyte scale in total.

自己紹介

いま何してる人？

★ 2023年4月～



の AI Lab. の Algorithms チーム（リサーチャー）

★ 2023年9月～ NII（国立情報学研究所）で特別研究員

★ 2023年～ 東工大で非常勤講師

自己紹介

いま何してる人？

★ 2023年4月～



の AI Lab. の Algorithms チーム（リサーチャー）

★ 2023年9月～NII（国立情報学研究所）で特別研究員

★ 2023年～東工大で非常勤講師

筑波大学との関わりは？

★ 2008年に情報学群情報科学類に入学しました

★ 2018年に博士（工学）で卒業しました

★ いまはもうない「記号計算研究室（SCORE 研）」の出身です

★ 学部1年の時に15単位かなんかしら取らなくて留年してます

★ WORD 編集部という、極めて良い組織、に所属していました

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 $\mathbb{F}[2]$:

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

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

Easy because just componentwise XOR:

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

- ▶ 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 $\mathbb{F}[2]$:

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

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

Easy because just componentwise XOR:

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

- ▶ 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 $\mathbb{F}[2^8]$.

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists y. x \cdot y = 1 \pmod{p}$$

証明は？

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists y. x \cdot y = 1 \pmod{p}$$

証明は？ 背理法で

1. y が存在しないまたは複数存在すると

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists !y. x \cdot y = 1 \pmod{p}$$

証明は？ 背理法で

1. y が存在しないまたは複数存在すると
2. $x \cdot z = x \cdot z' \pmod{p}$ なる相異なる数 $z < z'$ が鳩の巣原理から存在

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists ! y. x \cdot y = 1 \pmod{p}$$

証明は？ 背理法で

1. y が存在しないまたは複数存在すると
2. $x \cdot z = x \cdot z' \pmod{p}$ なる相異なる数 $z < z'$ が鳩の巣原理から存在
3. $x \cdot (z' - z) = 0 \pmod{p}$ なので x か $z' - z$ が p の倍数

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists y. x \cdot y = 1 \pmod{p}$$

証明は？ 背理法で

1. y が存在しないまたは複数存在すると
2. $x \cdot z = x \cdot z' \pmod{p}$ なる相異なる数 $z < z'$ が鳩の巣原理から存在
3. $x \cdot (z' - z) = 0 \pmod{p}$ なので x か $z' - z$ が p の倍数
4. しかし x も $z' - z$ も p 未満である。よって矛盾。

もうちょっと構成を詳しく説明します

まず簡単なところから。素数 p について $F[p]$ を構成するには？

- ▶ 足し算は $x + y \pmod{p}$ とすれば良い。加算逆元 $-x$ も簡単。
- ▶ 掛け算も $x \cdot y \pmod{p}$ としたい。乗算逆元 x^{-1} の保証は？
- ▶ 次の性質を使います：

$$\forall x. \exists y. x \cdot y = 1 \pmod{p}$$

証明は？ 背理法で

1. y が存在しないまたは複数存在すると
2. $x \cdot z = x \cdot z' \pmod{p}$ なる相異なる数 $z < z'$ が鳩の巣原理から存在
3. $x \cdot (z' - z) = 0 \pmod{p}$ なので x か $z' - z$ が p の倍数
4. しかし x も $z' - z$ も p 未満である。よって矛盾。

- ▶ もちろんフェルマーの小定理を用いても良いです：

$$\forall a. (a \pmod{p} \neq 0) \implies a^{p-1} = 1 \pmod{p}$$

有限体を作つてみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

有限体を作つてみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

有限体を作つてみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

これって、bit XOR やん！！（そうです）

\otimes	0	1
0	0	0
1	0	1

有限体を作ってみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

これって、bit XOR やん！！（そうです）

\otimes	0	1
0	0	0
1	0	1

これって、bit AND やん！！（そうです）

チェック項目:

- ① 加算逆元はある？ x に対して $-x$ があるか？

有限体を作ってみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

これって、bit XOR やん！！（そうです）

\otimes	0	1
0	0	0
1	0	1

これって、bit AND やん！！（そうです）

チェック項目:

- ① 加算逆元はある？ x に対して $-x$ があるか？ ある $x \oplus x = 0$

有限体を作ってみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

これって、bit XOR やん！！（そうです）

\otimes	0	1
0	0	0
1	0	1

これって、bit AND やん！！（そうです）

チェック項目:

- ① 加算逆元はある？ x に対して $-x$ があるか？ ある $x \oplus x = 0$
- ② 積逆元はある？ $x \neq 0$ に対して $\frac{1}{x}$ があるか？

有限体を作ってみよう $\mathbb{F}[2] = \{0, 1\}$

要素数が素数 p の場合は簡単

加算 $x \oplus y$ の定義 = 普通に足して剰余とる $(x + y) \bmod p$

乗算 $x \otimes y$ の定義 = 普通に掛けて剰余とる $(x \cdot y) \bmod p$

演算表

\oplus	0	1
0	0	1
1	1	0

これって、bit XOR やん！！（そうです）

\otimes	0	1
0	0	0
1	0	1

これって、bit AND やん！！（そうです）

チェック項目:

- ① 加算逆元はある？ x に対して $-x$ があるか？ ある $x \oplus x = 0$
- ② 積逆元はある？ $x \neq 0$ に対して $\frac{1}{x}$ があるか？ ある $1 \otimes 1 = 1$

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

$$0 \oplus x_0 = 0, \quad 1 \oplus x_1 = 0, \quad 2 \oplus x_2 = 0,$$

$$3 \oplus x_3 = 0, \quad 4 \oplus x_4 = 0, \quad 5 \oplus x_5 = 0,$$

$$6 \oplus x_6 = 0,$$

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

$$0 \oplus x_0 = 0, \quad 1 \oplus x_1 = 0, \quad 2 \oplus x_2 = 0,$$

$$3 \oplus x_3 = 0, \quad 4 \oplus x_4 = 0, \quad 5 \oplus x_5 = 0,$$

$$6 \oplus x_6 = 0,$$

掛け算 $x \otimes y$ は普通の掛け算をして 7 で割る: $(x \times y) \bmod 7$

全ての非零数 i に積逆元 x_i (つまり $\frac{1}{i}$) が存在するか調べよ:

$$1 \otimes x_1 = 1, \quad 2 \otimes x_2 = 1,$$

$$3 \otimes x_3 = 1, \quad 4 \otimes x_4 = 1, \quad 5 \otimes x_5 = 1,$$

$$6 \otimes x_6 = 1,$$

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

$$0 \oplus x_0 = 0, \quad 1 \oplus x_1 = 0, \quad 2 \oplus x_2 = 0,$$

$$3 \oplus x_3 = 0, \quad 4 \oplus x_4 = 0, \quad 5 \oplus x_5 = 0,$$

$$6 \oplus x_6 = 0,$$

掛け算 $x \otimes y$ は普通の掛け算をして 7 で割る: $(x \times y) \bmod 7$

全ての非零数 i に積逆元 x_i (つまり $\frac{1}{i}$) が存在するか調べよ:

$$1 \otimes x_1 = 1, \quad 2 \otimes x_2 = 1,$$

$$3 \otimes x_3 = 1, \quad 4 \otimes x_4 = 1, \quad 5 \otimes x_5 = 1,$$

$$6 \otimes x_6 = 1,$$

次を満たす数 α を (積の) 生成元と呼ぶ。これを求めよ

$\alpha, \alpha^2, \alpha^3, \alpha^4, \alpha^5, \alpha^6$ が全て異なる

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

答え

$$0 \oplus 0 = 0, 1 \oplus 6 = 0, 2 \oplus 5 = 0,$$

$$3 \oplus 4 = 0, 4 \oplus 3 = 0, 5 \oplus 2 = 0, 6 \oplus 1 = 0$$

$$\text{つまり } -1 = 6, -2 = 5, -3 = 4, -4 = 3, -5 = 2, -6 = 1$$

掛け算 $x \otimes y$ は普通の掛け算をして 7 で割る: $(x \times y) \bmod 7$

全ての非零数 i に積逆元 x_i (つまり $\frac{1}{i}$) が存在するか調べよ:

$$1 \otimes x_1 = 1, 2 \otimes x_2 = 1,$$

$$3 \otimes x_3 = 1, 4 \otimes x_4 = 1, 5 \otimes x_5 = 1,$$

$$6 \otimes x_6 = 1,$$

次を満たす数 α を (積の) 生成元と呼ぶ。これを求めよ

$$\alpha, \alpha^2, \alpha^3, \alpha^4, \alpha^5, \alpha^6 \text{ が全て異なる}$$

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

答え

$$0 \oplus 0 = 0, 1 \oplus 6 = 0, 2 \oplus 5 = 0,$$

$$3 \oplus 4 = 0, 4 \oplus 3 = 0, 5 \oplus 2 = 0, 6 \oplus 1 = 0$$

$$\text{つまり } -1 = 6, -2 = 5, -3 = 4, -4 = 3, -5 = 2, -6 = 1$$

掛け算 $x \otimes y$ は普通の掛け算をして 7 で割る: $(x \times y) \bmod 7$

全ての非零数 i に積逆元 x_i (つまり $\frac{1}{i}$) が存在するか調べよ:

答え

$$1 \otimes 1 = 1, 2 \otimes 4 = 1, 3 \otimes 5 = 1,$$

$$4 \otimes 2 = 1, 5 \otimes 3 = 1, 6 \otimes 6 = 1$$

$$\text{つまり } \frac{1}{1} = 1, \frac{1}{2} = 4, \frac{1}{3} = 5, \frac{1}{4} = 2, \frac{1}{5} = 3, \frac{1}{6} = 6$$

次を満たす数 α を (積の) 生成元と呼ぶ。これを求めよ

$\alpha, \alpha^2, \alpha^3, \alpha^4, \alpha^5, \alpha^6$ が全て異なる

練習問題: $\mathbb{F}[7] = \{0, 1, 2, 3, 4, 5, 6\}$ を構成しよう

足し算 $x \oplus y$ は普通の足し算をして 7 で割る: $(x + y) \bmod 7$

全ての数 i に加算逆元 x_i (つまり $-i$) は存在するか調べよ:

答え

$$0 \oplus 0 = 0, 1 \oplus 6 = 0, 2 \oplus 5 = 0,$$

$$3 \oplus 4 = 0, 4 \oplus 3 = 0, 5 \oplus 2 = 0, 6 \oplus 1 = 0$$

$$\text{つまり } -1 = 6, -2 = 5, -3 = 4, -4 = 3, -5 = 2, -6 = 1$$

掛け算 $x \otimes y$ は普通の掛け算をして 7 で割る: $(x \times y) \bmod 7$

全ての非零数 i に積逆元 x_i (つまり $\frac{1}{i}$) が存在するか調べよ:

答え

$$1 \otimes 1 = 1, 2 \otimes 4 = 1, 3 \otimes 5 = 1,$$

$$4 \otimes 2 = 1, 5 \otimes 3 = 1, 6 \otimes 6 = 1$$

$$\text{つまり } \frac{1}{1} = 1, \frac{1}{2} = 4, \frac{1}{3} = 5, \frac{1}{4} = 2, \frac{1}{5} = 3, \frac{1}{6} = 6$$

次を満たす数 α を (積の) 生成元と呼ぶ。これを求めよ

答え: $\alpha = 3; 3^2 = 2, 3^3 = 6, 3^4 = 4, 3^5 = 5, 3^6 = 1$ ($\alpha = 5$ も解)

もうちょっと構成を詳しく説明します

$p = 2$ の場合について確認: $\mathbb{F}[2] = \{0, 1\}$ です。

- ▶ 加算は $x + y \pmod{2}$ で、これは実質 bit-XOR
- ▶ 乗算は $x \cdot y \pmod{2}$ で、これは実質 bit-AND

$\mathbb{F}[2^8] = \{0, 1, 2, \dots, 255\}$ の場合はどうか?

もうちょっと構成を詳しく説明します

$p = 2$ の場合について確認: $\mathbb{F}[2] = \{0, 1\}$ です。

- ▶ 加算は $x + y \pmod{2}$ で、これは実質 bit-XOR
- ▶ 乗算は $x \cdot y \pmod{2}$ で、これは実質 bit-AND

$\mathbb{F}[2^8] = \{0, 1, 2, \dots, 255\}$ の場合はどうか?

- ▶ $x \cdot y \pmod{256}$ で乗算を定義すると、一般に逆元が存在しないです。
- ▶ 例: $2 \cdot y = 1 \pmod{256}$ とする y がない (2 の積逆元がない)。

もうちょっと構成を詳しく説明します

$p = 2$ の場合について確認: $\mathbb{F}[2] = \{0, 1\}$ です。

- ▶ 加算は $x + y \pmod{2}$ で、これは実質 bit-XOR
- ▶ 乗算は $x \cdot y \pmod{2}$ で、これは実質 bit-AND

$\mathbb{F}[2^8] = \{0, 1, 2, \dots, 255\}$ の場合はどうか?

- ▶ $x \cdot y \pmod{256}$ で乗算を定義すると、一般に逆元が存在しないです。
- ▶ 例: $2 \cdot y = 1 \pmod{256}$ とする y がない (2 の積逆元がない)。

この問題を克服するために

$$k_7x^7 + k_6x^6 + \cdots + k_1x + k_0 \quad (k_i \in \mathbb{F}[2])$$

という多項式全体 (濃度は 2^8) を使えるというのが大きいです。

もうちょっと構成を詳しく説明します

$p = 2$ の場合について確認: $\mathbb{F}[2] = \{0, 1\}$ です。

- ▶ 加算は $x + y \pmod{2}$ で、これは実質 bit-XOR
- ▶ 乗算は $x \cdot y \pmod{2}$ で、これは実質 bit-AND

$\mathbb{F}[2^8] = \{0, 1, 2, \dots, 255\}$ の場合はどうか?

- ▶ $x \cdot y \pmod{256}$ で乗算を定義すると、一般に逆元が存在しないです。
- ▶ 例: $2 \cdot y = 1 \pmod{256}$ とする y がない (2 の積逆元がない)。

この問題を克服するために

$$k_7x^7 + k_6x^6 + \cdots + k_1x + k_0 \quad (k_i \in \mathbb{F}[2])$$

という多項式全体 (濃度は 2^8) を使えるというのが大きいです。

7次多項式の積は7次を超えててしまうので、

素数っぽい振る舞いをする8次多項式 (=既約多項式) で割ります。

原始多項式というフェルマーの小定理っぽいのを満たすものもあります。

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

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

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

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

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

- ▶ 1-byte and 8-bits are isomorphic: $x \in \mathbb{F}[2^8] \cong \tilde{x} \in \begin{matrix} 1 \\ 8 \end{matrix} \mathbb{F}[2]$.
- ▶ There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \begin{matrix} 8 \\ \mathbb{F}[2] \end{matrix}$ i.e.,
$$\forall x, y \in \mathbb{F}[2^8]. \quad \left\{ \begin{array}{l} x + y = \mathcal{B}^{-1}(\mathcal{B}(x) + \mathcal{B}(y)), \\ x \cdot y = \mathcal{B}^{-1}(\mathcal{B}(x) \times \mathcal{B}(y)) \end{array} \right.$$

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 $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \begin{matrix} 1 \\ \boxed{\mathbb{F}[2]} \\ 8 \end{matrix}$ 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

$$\mathcal{B}(\mathcal{W}^{-1}) \stackrel{\mathbb{F}[2]}{\times} (\mathcal{B}(\mathcal{W}) \stackrel{\mathbb{F}[2]}{\times} \widetilde{D}) = \mathcal{B}(\mathcal{W}^{-1} \times \mathcal{W}) \times \widetilde{D} = \widetilde{D}.$$

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

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

$\begin{matrix} 1 \\ \boxed{\mathbb{F}[2]} \\ 8 \end{matrix}$

- ▶ There is an injective ring homomorphism $\mathcal{B} : \mathbb{F}[2^8] \rightarrow \mathbb{F}[2]^8$

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

⇓

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

\oplus is byte-array XOR.

XOR-based EC:

- ▶ 1-byte and 8-bits

- ▶ There is an injecti

演算の種類が一種類になった！
 しかもバイト列のXORするだけ
 これは SIMD化 が効きそう！

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

↓

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

\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	$\text{RS}(10, 4)$ by $\mathbb{F}[2^8]$	$\text{RS}(10, 4)$ by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \begin{matrix} 10 \\ \mathbb{F}[2^8] \end{matrix}$	$\mathcal{B}(\mathcal{V}) : 112 \begin{matrix} 80 \\ \mathbb{F}[2] \end{matrix}$
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	RS(10, 4) by $\mathbb{F}[2^8]$	RS(10, 4) by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \begin{matrix} 10 \\ \mathbb{F}[2^8] \end{matrix}$	$\mathcal{B}(\mathcal{V}) : 112 \begin{matrix} 80 \\ \mathbb{F}[2] \end{matrix}$
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-l>

♠ 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	RS(10, 4) by $\mathbb{F}[2^8]$	RS(10, 4) by $\mathbb{F}[2]$
Number of Core Operation	$\mathcal{V} : 14 \begin{matrix} 10 \\ \mathbb{F}[2^8] \end{matrix}$	$\mathcal{B}(\mathcal{V}) : 112 \begin{matrix} 80 \\ \mathbb{F}[2] \end{matrix}$
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-l>

♦ 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):

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

$$\Downarrow$$
$$\frac{P(a, b, c, d)}{v_1 \leftarrow a \oplus b; \\ v_2 \leftarrow a \oplus b \oplus c; \\ v_3 \leftarrow b \oplus c \oplus d; \\ \text{return}(v_1, v_2, v_3)}$$

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

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

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

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

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

$$\frac{P(a, b, c, d)}{v_1 \leftarrow a \oplus b; \quad v_2 \leftarrow a \oplus b \oplus c; \quad v_3 \leftarrow b \oplus c \oplus d; \quad \text{return}(v_1, v_2, v_3)}$$

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

- ★ "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):

$$\left(\begin{array}{cccc} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{array} \right) \times_{\mathbb{F}[2]} \begin{pmatrix} \vec{a} \\ \vec{b} \\ \vec{c} \\ \vec{d} \end{pmatrix} = \begin{pmatrix} \vec{a} \oplus \vec{b} \\ \vec{a} \oplus \vec{b} \oplus \vec{c} \\ \vec{b} \oplus \vec{c} \oplus \vec{d} \end{pmatrix}$$

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

- ★ "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.

SLP (=上から下に実行するだけ言語) の最適化の例

プログラムの意味を変えず、演算回数を減らす:

```
x1 ← a ⊕ b ⊕ c ⊕ d ⊕ x;  
x2 ← b ⊕ c ⊕ a ⊕ d ⊕ y;  
x3 ← a ⊕ c ⊕ z ⊕ b;  
return(x1, x2, x3); # 11 個
```

SLP (=上から下に実行するだけ言語) の最適化の例

プログラムの意味を変えず、演算回数を減らす:

```
x1 ← a ⊕ b ⊕ c ⊕ d ⊕ x;  
x2 ← b ⊕ c ⊕ a ⊕ d ⊕ y;  
x3 ← a ⊕ c ⊕ z ⊕ b;  
return(x1, x2, x3); # 11 個
```

変数は追加して良い:

```
t1 ← a ⊕ b ⊕ c;  
x1 ← t1 ⊕ d ⊕ x;  
x2 ← t1 ⊕ d ⊕ y;  
x3 ← t1 ⊕ z;  
return(x1, x2, x3); # 7 個
```

SLP (=上から下に実行するだけ言語) の最適化の例

プログラムの意味を変えず、演算回数を減らす:

```
x1 ← a ⊕ b ⊕ c ⊕ d ⊕ x;  
x2 ← b ⊕ c ⊕ a ⊕ d ⊕ y;  
x3 ← a ⊕ c ⊕ z ⊕ b;  
return(x1, x2, x3); # 11 個
```

変数は追加して良い:

```
t1 ← a ⊕ b ⊕ c;  
x1 ← t1 ⊕ d ⊕ x;  
x2 ← t1 ⊕ d ⊕ y;  
x3 ← t1 ⊕ z;  
return(x1, x2, x3); # 7 個
```

まだ削れる:

```
t1 ← a ⊕ b ⊕ c;  
t2 ← t1 ⊕ d;  
x1 ← t2 ⊕ x;  
x2 ← t2 ⊕ y;  
x3 ← t1 ⊕ z;  
return(x1, x2, x3); # 6 個
```

XOR Optimization: Reducing XORs

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

$$\frac{P \quad \#\oplus = 8}{\begin{aligned} v_1 &\leftarrow a \oplus b; \\ v_2 &\leftarrow a \oplus b \oplus c; \\ v_3 &\leftarrow a \oplus b \oplus c \oplus d; \\ v_4 &\leftarrow b \oplus c \oplus d; \end{aligned}} \qquad \qquad \frac{Q}{}$$

return(v_1, v_2, v_3, v_4)

XOR Optimization: Reducing XORs

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

$$\frac{P \quad \#\oplus = 8}{\begin{aligned} v_1 &\leftarrow a \oplus b; \\ v_2 &\leftarrow a \oplus b \oplus c; \\ v_3 &\leftarrow a \oplus b \oplus c \oplus d; \\ v_4 &\leftarrow b \oplus c \oplus d; \end{aligned}} \qquad \Rightarrow \qquad \frac{Q}{v_1 \leftarrow a \oplus b;}$$

return(v_1, v_2, v_3, v_4)

XOR Optimization: Reducing XORs

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

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

return(v_1, v_2, v_3, v_4)

XOR Optimization: Reducing XORs

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

$$\frac{P \quad \#\oplus = 8}{\begin{aligned} v_1 &\leftarrow a \oplus b; \\ v_2 &\leftarrow a \oplus b \oplus c; \\ v_3 &\leftarrow a \oplus b \oplus c \oplus d; \\ v_4 &\leftarrow b \oplus c \oplus d; \end{aligned}} \quad \Rightarrow \quad \frac{Q}{\begin{aligned} v_1 &\leftarrow a \oplus b; \\ v_2 &\leftarrow v_1 \oplus c; \\ v_3 &\leftarrow v_2 \oplus d; \end{aligned}}$$

return(v_1, v_2, v_3, v_4)

XOR Optimization: Reducing XORs

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

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

$$\frac{Q \quad \#\oplus = 4}{\begin{aligned} v_1 &\leftarrow a \oplus b; \\ v_2 &\leftarrow v_1 \oplus c; \\ v_3 &\leftarrow v_2 \oplus d; \\ v_4 &\leftarrow v_3 \oplus a; \end{aligned}} \quad \because (a \oplus b \oplus c \oplus d) \oplus a = b \oplus c \oplus d. \quad \text{return}(v_1, v_2, v_3, v_4)$$

XOR Optimization: Reducing XORs

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

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

- ▶ 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.

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

- ▶ 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.

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

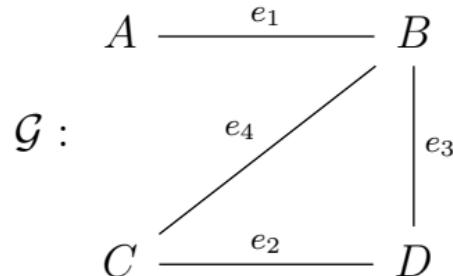
- ▶ P and Q are equivalent: $\llbracket P \rrbracket = \llbracket Q \rrbracket$.
- ▶ Intuitively, Q ($\#\oplus(Q) = 4$) runs faster than P ($\#\oplus(P) = 8$).

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 $\#\oplus(Q)$.*

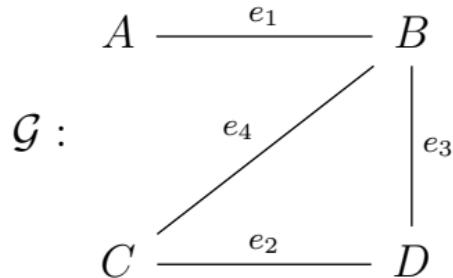
XOR最適化のNP完全性についてもうちょっと

Vertex Cover Problem という古典的なNP完全問題を使います。



XOR最適化のNP完全性についてもうちょっと

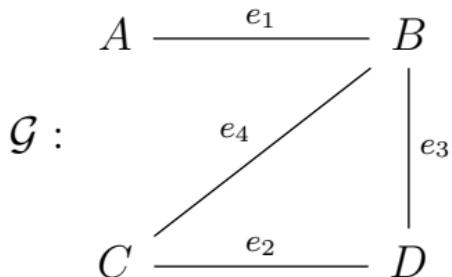
Vertex Cover Problem という古典的な NP 完全問題を使います。



このグラフ \mathcal{G} については、

- ▶ 頂点セット $\{B, C\}$ で、全ての辺をカバーできます
- ▶ 他の頂点セット $\{B, D\}$ でも、カバーできます

XOR最適化のNP完全性についてもうちょっと

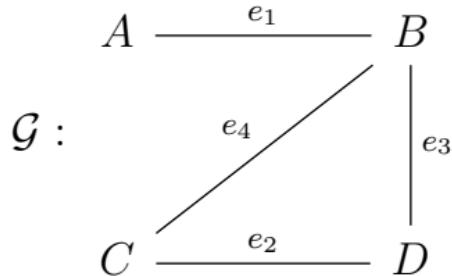


グラフ \mathcal{G} から、次のSLP $P_{\mathcal{G}}$ を作ります:

$$P_{\mathcal{G}} : \begin{aligned} e_1 &\leftarrow p \oplus A \oplus B; \\ e_2 &\leftarrow p \oplus C \oplus D; \\ e_3 &\leftarrow p \oplus B \oplus D; \\ e_4 &\leftarrow p \oplus B \oplus C; \end{aligned}$$

これをXOR最適化すると、最小頂点セットが実は現れます。

XOR最適化のNP完全性についてもうちょっと

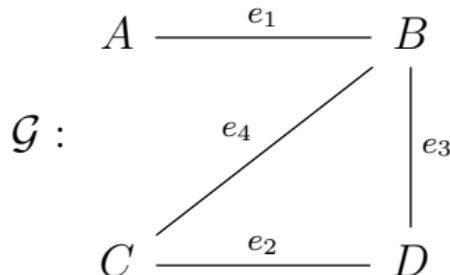


グラフ \mathcal{G} から、次の SLP $P_{\mathcal{G}}$ を作ります:

$$P_{\mathcal{G}} : \begin{array}{ll} e_1 \leftarrow p \oplus A \oplus B; & p_B \leftarrow p \oplus B; \\ e_2 \leftarrow p \oplus C \oplus D; & e_1 \leftarrow p_B \oplus A; \\ e_3 \leftarrow p \oplus B \oplus D; & e_3 \leftarrow p_B \oplus D; \\ e_4 \leftarrow p \oplus B \oplus C; & e_4 \leftarrow p_B \oplus C; \\ & p_C \leftarrow p \oplus C; \\ & e_2 \leftarrow p_C \oplus D; \end{array} \Rightarrow P^{\text{opt}} :$$

これを XOR 最適化すると、最小頂点セットが実は現れます。

XOR最適化のNP完全性についてもうちょっと



グラフ \mathcal{G} から、次の SLP $P_{\mathcal{G}}$ を作ります:

$$P_{\mathcal{G}} : \begin{array}{ll} e_1 \leftarrow p \oplus A \oplus B; & p_B \leftarrow p \oplus B; \\ e_2 \leftarrow p \oplus C \oplus D; & e_1 \leftarrow p_B \oplus A; \\ e_3 \leftarrow p \oplus B \oplus D; & e_3 \leftarrow p_B \oplus D; \\ e_4 \leftarrow p \oplus B \oplus C; & e_4 \leftarrow p_B \oplus C; \\ & p_C \leftarrow p \oplus C; \\ & e_2 \leftarrow p_C \oplus D; \end{array} \Rightarrow P^{\text{opt}} :$$

これを XOR 最適化すると、最小頂点セットが実は現れます。 実際に示しているのは、いつでも右のような形にできる、という正規化補題です。

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:

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

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

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

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:

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

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

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

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

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

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

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:

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

The commutative version of REPAIR accommodates

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

In the paper, we extend it to XORREPAIR by accommodating

$$\text{Cancellativity: } x \oplus x \oplus y = y.$$

文法圧縮との出会い▷ Tozawa & Minamide, FOSSACS'07.

https://link.springer.com/chapter/10.1007%2F978-3-540-71389-0_25

Complexity Results on Balanced Context-Free Languages

Akihiko Tozawa¹ and Yasuhiko Minamide²

¹ IBM Research,
Tokyo Research Laboratory, IBM Japan, Ltd.

² Department of Computer Science
University of Tsukuba

Abstract. Some decision problems related to balanced context-free languages are important for their application to the static analysis of programs generating XML strings. One such problem is the balancedness problem which decides whether or not the language of a given context-free grammar (CFG) over a paired alphabet is balanced. Another important problem is the validation problem which decides whether or not the language of a CFG is contained by that of a regular hedge grammar (RHG). This paper gives two new results; (1) the balancedness problem is in PTIME; and (2) the CFG-RHG containment problem is 2EXPTIME-complete.

文法圧縮との出会い▷ Tozawa & Minamide, FOSSACS'07.

https://link.springer.com/chapter/10.1007%2F978-3-540-71389-0_25

Complexity Results on Balanced Context-Free Languages

論文の前半では:

Akihiko Tozawa¹ and Yasuhiko Minamide²

¹ IBM Research,
Tokyo Research Laboratory, IBM Japan, Ltd.

² Department of Computer Science
University of Tsukuba

- ▶ CFG $G \subseteq?$ Dyck を解く
- ▶ ただ解くだけでなく PTIME で解くために文法圧縮の技を使っている

Abstract. Some decision problems related to balanced context-free languages are important for their application to the static analysis of programs generating XML strings. One such problem is the balancedness problem which decides whether or not the language of a given context-free grammar (CFG) over a paired alphabet is balanced. Another important problem is the validation problem which decides whether or not the language of a CFG is contained by that of a regular hedge grammar (RHG). This paper gives two new results; (1) the balancedness problem is in PTIME; and (2) the CFG-RHG containment problem is 2EXPTIME-complete.

文法圧縮との出会い▷ Tozawa & Minamide, FOSSACS'07.

https://link.springer.com/chapter/10.1007%2F978-3-540-71389-0_25

Complexity Results on Balanced Context-Free Languages

論文の前半では:

Akihiko Tozawa¹ and Yasuhiko Minamide²

¹ IBM Research,
Tokyo Research Laboratory, IBM Japan, Ltd.

² Department of Computer Science
University of Tsukuba

- ▶ CFG $G \subseteq?$ Dyck を解く
- ▶ ただ解くだけでなく PTIME で解くために文法圧縮の技を使っている

論文の後半: CFG $\subseteq?$ RHG

Abstract. Some decision problems related to balanced context-free languages are important for their application to the static analysis of programs generating XML strings. One such problem is the balancedness problem which decides whether or not the language of a given context-free grammar (CFG) over a paired alphabet is balanced. Another important problem is the validation problem which decides whether or not the language of a CFG is contained by that of a regular hedge grammar (RHG). This paper gives two new results; (1) the balancedness problem is in PTIME; and (2) the CFG-RHG containment problem is 2EXPTIME-complete.

文法圧縮との出会い▷ Tozawa & Minamide, FOSSACS'07.

https://link.springer.com/chapter/10.1007%2F978-3-540-71389-0_25

Complexity Results on Balanced Context-Free Languages

論文の前半では:

Akihiko Tozawa¹ and Yasuhiko Minamide²

¹ IBM Research,
Tokyo Research Laboratory, IBM Japan, Ltd.
² Department of Computer Science
University of Tsukuba

- ▶ CFG $G \subseteq?$ Dyck を解く
- ▶ ただ解くだけでなく PTIME で解くために文法圧縮の技を使っている

Abstract. Some decision problems related to balanced context-free languages are important for their application to the static analysis of programs generating XML strings. One such problem is the balancedness problem which decides whether or not the language of a given context-free grammar (CFG) over a paired alphabet is balanced. Another important problem is the validation problem which decides whether or not the language of a CFG is contained by that of a regular hedge grammar (RHG). This paper gives two new results; (1) the balancedness problem is in PTIME; and (2) the CFG-RHG containment problem is 2EXPTIME-complete.

論文の後半: CFG $\subseteq?$ RHG
Uezato & Minamide DLT'16 で
CFG $\subseteq?$ Superdeterministic PDA に
拡張済み

文法圧縮との出会い▷ Tozawa & Minamide, FOSSACS'07.

https://link.springer.com/chapter/10.1007%2F978-3-540-71389-0_25

Complexity Results on Balanced Context-Free Languages

論文の前半では:

Akihiko Tozawa¹ and Yasuhiko Minamide²

¹ IBM Research,
Tokyo Research Laboratory, IBM Japan, Ltd.
² Department of Computer Science
University of Tsukuba

- ▶ CFG $G \subseteq?$ Dyck を解く
- ▶ ただ解くだけでなく PTIME で解くために文法圧縮の技を使っている

論文の後半: CFG $\subseteq?$ RHG

Uezato & Minamide DLT'16 で
CFG $\subseteq?$ Superdeterministic PDA に
拡張済み

Abstract. Some decision problems related to balanced context-free languages are important for their application to the static analysis of programs generating XML strings. One such problem is the balancedness problem which decides whether or not the language of a given context-free grammar (CFG) over a paired alphabet is balanced. Another important problem is the validation problem which decides whether or not the language of a CFG is contained by that of a regular hedge grammar (RHG). This paper gives two new results; (1) the balancedness problem is in PTIME; and (2) the CFG-RHG containment problem is 2EXPTIME-complete.

文法圧縮そのものについては次がオススメ:

The smallest grammar problem, Charikar+, 2005

<https://ieeexplore.ieee.org/document/1459058>

Memory Access Optimization: MultiSLP

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

Quiz: How many times will this program access memory?

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

Memory Access Optimization: MultiSLP

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

Quiz: How many times will this program access memory?

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

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

Quiz: How many times will this program access memory?

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

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

Quiz: How many times will this program access memory?

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

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 $\#_{\text{mem}} = 5$.

```
⊕4(A, B, C, D: [byte]) {
    var v = Array::new(A.len);
    for i in [0..A.len):
        byte r = A[i] ^ B[i]
        r = r ^ C[i];
        v[i] = r ^ 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. $\#_{\text{mem}}$?

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

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. $\#_{\text{mem}}$?

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

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. $\#_{\text{mem}}$?

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

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 $\#_{\text{mem}}(Q)$.*

メモ: Deforestation

$$\frac{P}{}$$

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

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

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

$$\frac{P'}{t \leftarrow a \oplus b \oplus c \oplus d;}$$

$$v_1 \leftarrow t \oplus e;$$

$$v_2 \leftarrow t \oplus f;$$

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

$$\frac{Q}{t \leftarrow \oplus_4(a, b, c, d);}$$

$$v_1 \leftarrow t \oplus e;$$

$$v_2 \leftarrow t \oplus f;$$

$$\#_{\text{mem}}(Q) = 11$$

$P' \Rightarrow Q$ でやっている合成による最適化（中間データの削除）は
関数プログラミングでは「Deforestation」と呼ばれる。

メモ: Deforestation

$$\frac{P}{\begin{array}{l} v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \#_{\text{mem}}(P) = 24 \end{array}} \quad \Rightarrow \quad \frac{P'}{\begin{array}{l} t \leftarrow a \oplus b \oplus c \oplus d; \\ v_1 \leftarrow t \oplus e; \\ v_2 \leftarrow t \oplus f; \\ \#_{\text{mem}}(P') = 15 \end{array}} \quad \Rightarrow \quad \frac{Q}{\begin{array}{l} t \leftarrow \oplus_4(a, b, c, d); \\ v_1 \leftarrow t \oplus e; \\ v_2 \leftarrow t \oplus f; \\ \#_{\text{mem}}(Q) = 11 \end{array}}$$

$P' \Rightarrow Q$ でやっている合成による最適化（中間データの削除）は
関数プログラミングでは「Deforestation」と呼ばれる。

DEFORESTATION: TRANSFORMING PROGRAMS TO ELIMINATE TREES *

Philip WADLER

Department of Computer Science, University of Glasgow, Glasgow G12 8QQ, UK

Abstract. An algorithm that transforms programs to eliminate intermediate trees is presented. The algorithm applies to any term containing only functions with definitions in a given syntactic form, and is suitable for incorporation in an optimizing compiler.

メモ: Deforestation

P

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

P'

$$\begin{aligned} t &\leftarrow a \oplus b \oplus c \oplus d; \\ v_1 &\leftarrow t \oplus e; \\ v_2 &\leftarrow t \oplus f; \\ \#_{\text{mem}}(P') &= 15 \end{aligned}$$

Q

$$\begin{aligned} t &\leftarrow \oplus_4(a, b, c, d); \\ v_1 &\leftarrow t \oplus e; \\ v_2 &\leftarrow t \oplus f; \\ \#_{\text{mem}}(Q) &= 11 \end{aligned}$$

$P' \Rightarrow Q$ でやっている合成による最適化（中間データの削除）は
関数プログラミングでは「Deforestation」と呼ばれる。

sum (map (\x. x * x) (upto 1 n)) \Rightarrow

h 0 1 n

where

h a m n = if m > n

then a

else h(a + square m)(m + 1)n.

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

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

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

★ α appears once in the program

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

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

⋮

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

☆ α appears once in the program

Example.

$$v_1 \leftarrow a \oplus b \oplus c \oplus d \oplus e; \\ v_2 \leftarrow a \oplus b \oplus c \oplus d \oplus f; \\ \#_{\text{mem}}(24)$$

$\xrightarrow{\text{REPAIR}}$

$$t_1 \leftarrow a \oplus b; \\ t_2 \leftarrow t_1 \oplus c; \\ t_3 \leftarrow t_2 \oplus d; \\ v_1 \leftarrow t_3 \oplus e; \\ v_2 \leftarrow t_3 \oplus f; \\ \#_{\text{mem}}(15)$$

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

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

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

⋮

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

☆ α appears once in the program

Example.

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

$$\begin{aligned} t_1 &\leftarrow a \oplus b; \\ t_2 &\leftarrow t_1 \oplus c; \\ t_3 &\leftarrow t_2 \oplus d; \\ v_1 &\leftarrow t_3 \oplus e; \\ v_2 &\leftarrow t_3 \oplus f; \end{aligned} \quad \#_{\text{mem}}(15)$$

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

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

Our Heuristic: XOR Fusion

We fuse XORs when the following holds:

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

⋮

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

☆ α appears once in the program

Example.

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

$$\begin{aligned} t_1 &\leftarrow a \oplus b; \\ t_2 &\leftarrow t_1 \oplus c; \\ t_3 &\leftarrow t_2 \oplus d; \\ v_1 &\leftarrow t_3 \oplus e; \\ v_2 &\leftarrow t_3 \oplus f; \end{aligned} \quad \#_{\text{mem}}(15) \quad \xrightarrow{\text{fuse}(t_1)}$$

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

$$\begin{aligned} \text{fuse}(t_2) \Rightarrow & \quad t_3 \leftarrow \oplus_4(a, b, c, d); \\ & v_1 \leftarrow t_3 \oplus e; \\ & v_2 \leftarrow t_3 \oplus f; \end{aligned} \quad \#_{\text{mem}}(11) \quad \xrightarrow{\text{NOT fuse}(t_3) \text{ by } \star}$$

$$\begin{aligned} & v_1 \leftarrow \oplus_5(a, b, c, d, e); \\ & v_2 \leftarrow \oplus_5(a, b, c, d, f); \end{aligned} \quad \#_{\text{mem}}(12)$$

Cache Optimization: SLP + LRU Cache

Metric $\#_{\text{I/O}}(K, -)$: 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 :

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

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

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

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

return(v_2, v_3, v_4);

Cache Optimization: SLP + LRU Cache

Metric $\#_{I/O}(K, \cdot)$: 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 $\#_{\text{I/O}}(K, \cdot)$: 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 . $\#_{\text{I/O}} = 0$.
- ▶ $\mathcal{R}(x)$: Cache miss. Evict LRU to mem. and read x from mem. $\#_{\text{I/O}} = 2$.
- ▶ $\mathcal{W}(x)$: Cache miss. Evict LRU to mem. and write x to cache. $\#_{\text{I/O}} = 1$.

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

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

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

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

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

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 :

$$v_1 \leftarrow A \oplus B; \quad *_1 *_2 *_3 *_4 \xrightarrow[2]{\mathcal{R}(A)} *_2 *_3 *_4 A \xrightarrow[2]{\mathcal{R}(B)} *_3 *_4 AB \xrightarrow[1]{\mathcal{W}(v_1)} *$$

$$v_2 \leftarrow \oplus(E, D, A); \quad *_4 AB v_1$$

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

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

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 :

$$\begin{array}{lll} v_1 \leftarrow A \oplus B; & *_1 *_2 *_3 *_4 \xrightarrow[2]{\mathcal{R}(A)} *_2 *_3 *_4 A \xrightarrow[2]{\mathcal{R}(B)} *_3 *_4 AB \xrightarrow[1]{\mathcal{W}(v_1)} \\ v_2 \leftarrow \oplus(E, D, A); & *_4 ABv_1 \xrightarrow[2]{\mathcal{R}(E)} ABv_1E \xrightarrow[2]{\mathcal{R}(D)} Bv_1ED \xrightarrow[2]{\mathcal{R}(A)} v_1EDA \xrightarrow[1]{\mathcal{W}(v_2)} \\ v_3 \leftarrow v_1 \oplus E; & EDAv_2 \xrightarrow[2]{\mathcal{R}(v_1)} DA v_2 v_1 \xrightarrow[2]{\mathcal{R}(E)} Av_2 v_1 E \xrightarrow[1]{\mathcal{W}(v_3)} \\ v_4 \leftarrow v_1 \oplus C; & v_2 v_1 E v_3 \xrightarrow[0]{\mathcal{H}(v_1)} v_2 E v_3 v_1 \xrightarrow[2]{\mathcal{R}(C)} E v_3 v_1 C \xrightarrow[1]{\mathcal{W}(v_4)} \\ \text{return}(v_2, v_3, v_4); & v_3 v_1 C v_4 \implies \#_{I/O}(4, P) = 20. \end{array}$$

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*.

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

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	$v_1 \leftarrow A \oplus B;$
$v_2 \leftarrow \oplus(E, D, A);$	[7]	assignment	$v_2 \leftarrow \oplus(E, D, A);$
$v_3 \leftarrow v_1 \oplus E;$	[5]		$v_3 \leftarrow v_1 \oplus E;$
$v_4 \leftarrow v_1 \oplus C;$	[3]		$v_1 \leftarrow v_1 \oplus C;$
$\text{return}(v_2, v_3, v_4);$			$\text{return}(v_2, v_3, v_1);$

$$\begin{array}{c} \xrightarrow[0]{\mathcal{H}(v_1)} v_2Ev_3v_1 \xrightarrow[2]{\mathcal{R}(C)} Ev_3v_1C \xrightarrow[1]{\mathcal{W}(v_4)} v_3v_1Cv_4 \\ \Downarrow \\ \xrightarrow[0]{\mathcal{H}(v_1)} v_2Ev_3v_1 \xrightarrow[2]{\mathcal{R}(C)} Ev_3v_1C \xrightarrow[0]{\mathcal{H}(v_1)} Ev_3Cv_1 \end{array}$$

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.

	#I/O		#I/O
$v_1 \leftarrow A \oplus B;$	[5]	$v_2 \leftarrow \oplus(A, D, E);$	[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]
return (v_2, v_3, v_4);	20	return (v_2, v_3, v_4);	14

Reordering

Next Approach: Reordering Statements and Arguments

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

	#I/O		#I/O
$v_1 \leftarrow A \oplus B;$	[5]	$v_2 \leftarrow \oplus(A, D, E);$	[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]
return(v_2, v_3, v_4);	20	return(v_2, v_3, v_4);	14

Using *Pebble Game*, we can integrate  Recycling Variables and
Reordering

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

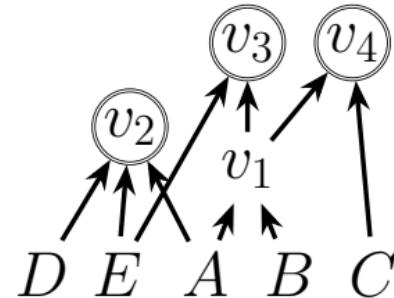
Next Approach: Reordering Statements and Arguments

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

	#I/O		#I/O
$v_1 \leftarrow A \oplus B;$	[5]	$v_2 \leftarrow \oplus(A, D, E);$	[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]
return(v_2, v_3, v_4);	20	return(v_2, v_3, v_4);	14

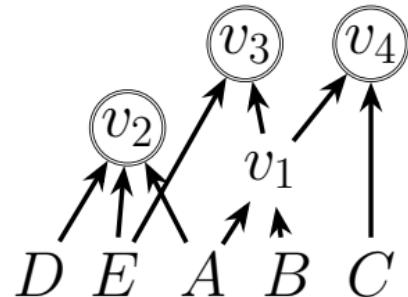
Using *Pebble Game*, we can integrate Recycling Variables and
Reordering

- ★ 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

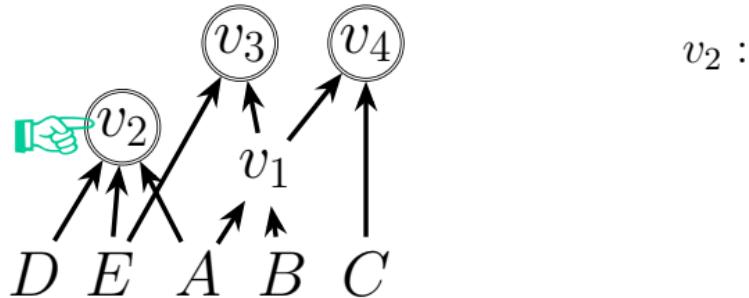
Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



Example: Evaluating strategy based on Depth-first-search

Pebble Game & Intractability of Optimization Problem

Playing Pebble Game = Deciding Evaluation Order + Variable Recycling



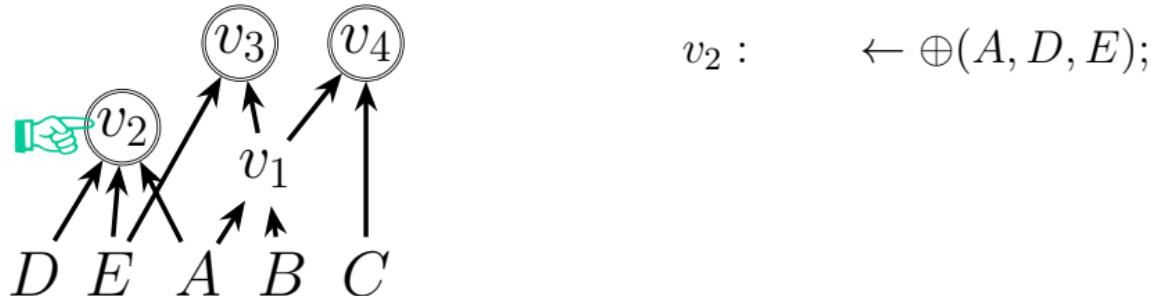
$v_2 :$

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



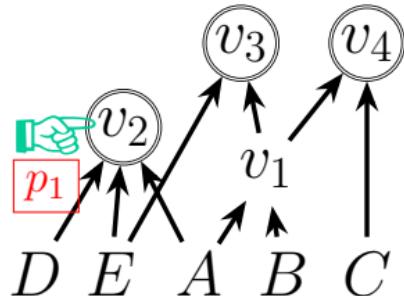
$v_2 : \quad \leftarrow \oplus(A, D, E);$

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



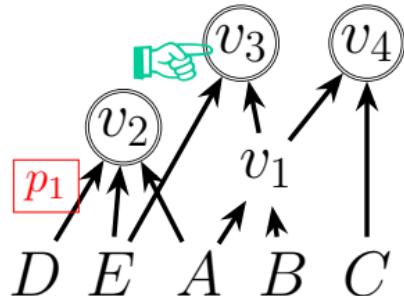
$$v_2 : \quad p_1 \leftarrow \oplus(A, D, E);$$

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



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

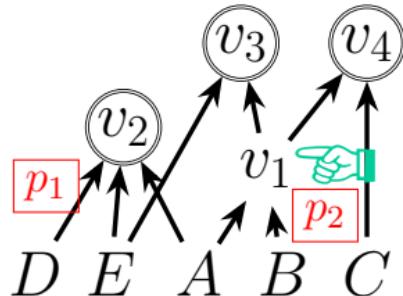
$$v_3 : \quad \leftarrow E \oplus$$

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



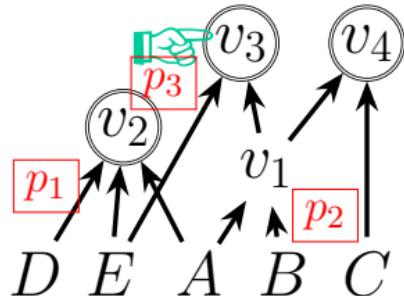
$$\begin{aligned} v_2 : \quad p_1 &\leftarrow \oplus(A, D, E); \\ v_1 : \quad p_2 &\leftarrow A \oplus B; \\ v_3 : \quad &\leftarrow E \oplus \end{aligned}$$

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



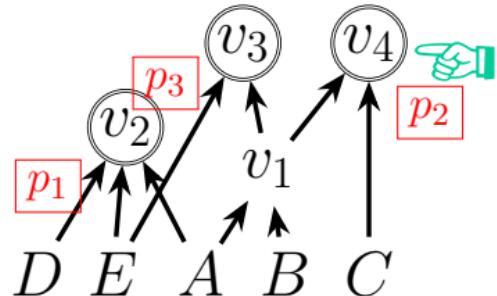
$$\begin{aligned}v_2 : \quad p_1 &\leftarrow \oplus(A, D, E); \\v_1 : \quad p_2 &\leftarrow A \oplus B; \\v_3 : \quad p_3 &\leftarrow E \oplus p_2;\end{aligned}$$

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



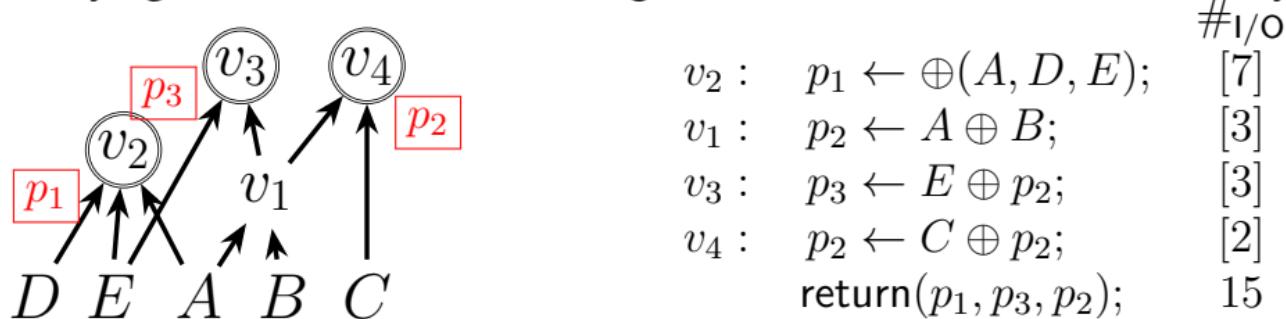
- $v_2 : p_1 \leftarrow \oplus(A, D, E);$
- $v_1 : p_2 \leftarrow A \oplus B;$
- $v_3 : p_3 \leftarrow E \oplus p_2;$
- $v_4 : p_2 \leftarrow C \oplus p_2;$

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
$v_2 : p_1 \leftarrow \oplus(A, D, E);$	[7]
$v_1 : p_2 \leftarrow A \oplus B;$	[3]
$v_3 : p_3 \leftarrow E \oplus p_2;$	[3]
$v_4 : p_2 \leftarrow C \oplus p_2;$	[2]
return(p_1, p_3, p_2);	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.

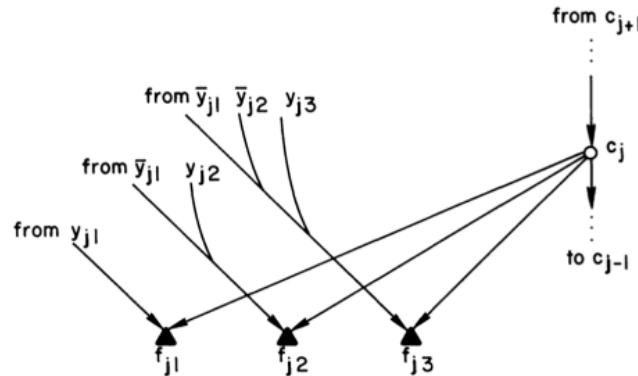
Pebble Game は 何から勉強すれば良い?

1975年に出版された

Complete Register Allocation Problems, Ravi Sethi

がオススメです。 <https://pubs.siam.org/doi/abs/10.1137/0204020>

こんな感じに図を書いて、3-SAT を pebble game 化します:



- direct descendants of final node (not shown)
- ▲ direct descendants of initial node (not shown)

FIG. 6. The subdag that checks if clause j is true

Pebble Game は 何から勉強すれば良い?

1975 年に出版された

Complete Register Allocation Problems, Ravi Sethi

がオススメです。<https://pubs.siam.org/doi/abs/10.1137/0204020>

これを I/O を考えるために拡張したのが

I/O complexity: The red-blue pebble game, Jia-Wei & Kung, STOC'81

<https://dl.acm.org/doi/10.1145/800076.802486>

Pebble Game は 何から勉強すれば良い?

1975 年に出版された

Complete Register Allocation Problems, Ravi Sethi

がオススメです。 <https://pubs.siam.org/doi/abs/10.1137/0204020>

これを I/O を考えるために拡張したのが

I/O complexity: The red-blue pebble game, Jia-Wei & Kung, STOC'81

<https://dl.acm.org/doi/10.1145/800076.802486>

よりモダンな Red-blue pebble game の話は

▶ *On the Hardness of Red-Blue Pebble Games*

Papp & Wattenhofer, SPAA'20

▶ *Red-blue pebbling revisited: near optimal parallel matrix-matrix multip.*

Kwasniewski+, SC'19

Pebble Game は 何から勉強すれば良い?

1975年に出版された

Complete Register Allocation Problems, Ravi Sethi

がオススメです。<https://pubs.siam.org/doi/abs/10.1137/0204020>

これを I/O を考えるために拡張したのが

I/O complexity: The red-blue pebble game, Jia-Wei & Kung, STOC'81

<https://dl.acm.org/doi/10.1145/800076.802486>

よりモダンな Red-blue pebble game の話は

▶ *On the Hardness of Red-Blue Pebble Games*

Papp & Wattenhofer, SPAA'20

▶ *Red-blue pebbling revisited: near optimal parallel matrix-matrix multip.*

Kwasniewski+, SC'19

教科書なら *Models Of Computation*, J. E. Savage

<http://cs.brown.edu/people/jsavage/book/>

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			
$\#_{\text{mem}}$	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

Base

\mathcal{B} -Byte Blocking for Cache Efficiency

```
v1 = xor(A, B);           for i ← 0 .. (A.len /  $\mathcal{B}$ ) {  
v2 = xor(v1, C, D);    v1[i] = xor(A[i], B[i]);  
return(v1, v2);          v2[i] = xor(v1[i], C[i], D[i]);  
}                         return(v1, v2);  
}
```

\mathcal{B} =

\mathcal{B} =

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

$\mathcal{B} = 2K$: #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			
$\#_{\text{mem}}$	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 + F	RePair + Fuse + ReBuild
# \oplus	755	Why smaller blocks are slower than the large one?		
#mem	2265	Pros: Smaller blocks, ▶ More cache-able blocks $\frac{32K}{\mathcal{B}}$.		
$\mathcal{B} = 512$: #I/O($K = 64$)	570	Cons: Smaller blocks, ▶ Due to cache conflicts, using cache identically is more difficult.		
Throughput (GB/s)	3.10	▶ Latency penalty becomes totally large.		
$\mathcal{B} = 1K$: #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 + Pebbling
# \oplus	755	385		
#mem	2265	1155		
$\mathcal{B} = 512 :$ #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 + 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		
#mem	2265	1155		
$\mathcal{B} = 512 :$	#I/O ($K = 64$)	570	1231	
	Throughput (GB/s)	3.10	4.18	
$\mathcal{B} = 1K :$	#I/O ($K = 32$)	1262	1465	
	Throughput (GB/s)	4.03	4.36	
$\mathcal{B} = 2K :$	#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	
#mem	2265	1155	677	
$\mathcal{B} = 512 :$	#I/O ($K = 64$)	570	1231	936
	Throughput (GB/s)	3.10	4.18	
$\mathcal{B} = 1K :$	#I/O ($K = 32$)	1262	1465	1086
	Throughput (GB/s)	4.03	4.36	
$\mathcal{B} = 2K :$	#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

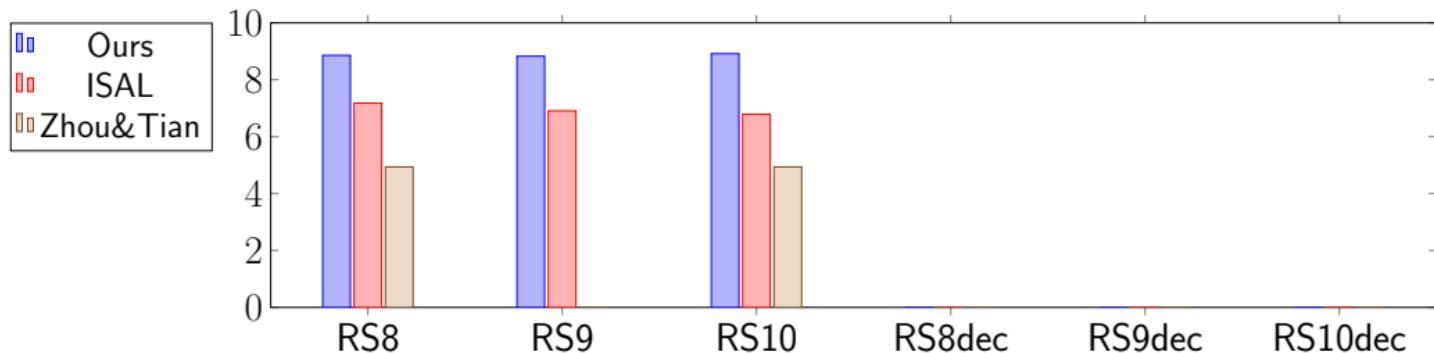
Improvements by heuristics for the encoding SLP on Intel PC

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
	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
				8.55

Throughput Comparison (Intel + 1K-Blocking)

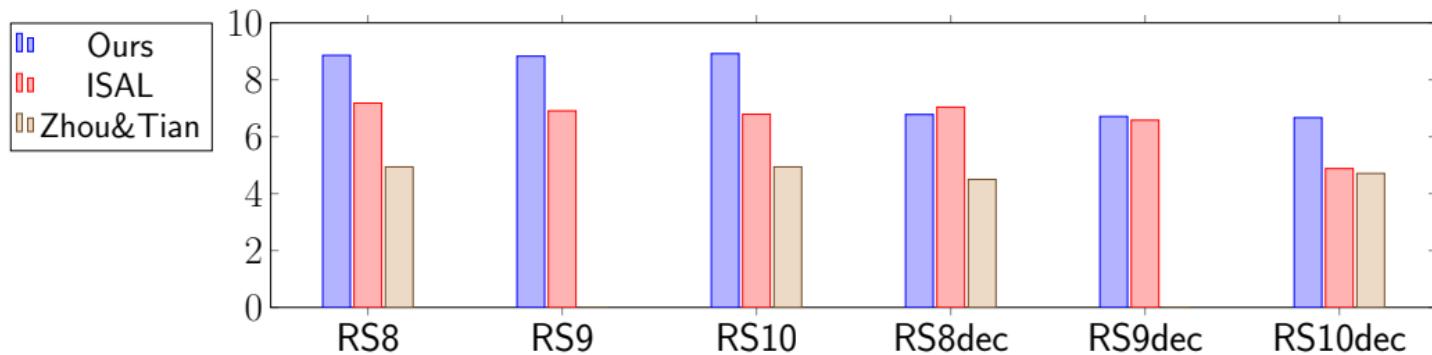
Enc	#mem	#I/O	Ours	ISA-L v2.30	Zhou & Tian
RS(8, 4)	543	585	8.86 GB/s	7.18 GB/s	4.94 GB/s
RS(9, 4)	611	671	8.83	6.91	N/A in their paper
RS(10, 4)	677	779	8.92	6.79	4.94



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.

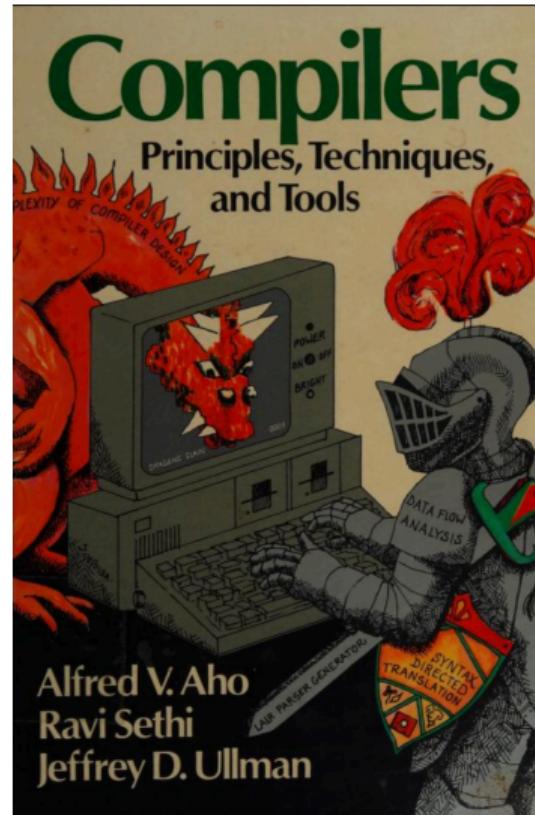
オススメの教科書: コンパイラ編

Aho & Ullman,

Principles of Compiler Design



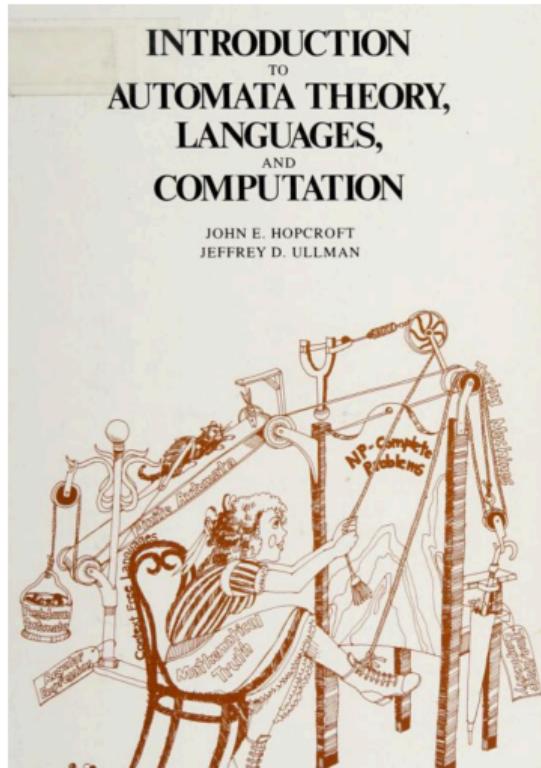
Aho & Sethi & Ullman,
Principles of Compiler Design



オススメの教科書: オートマトン & 計算量理論 (NP やらなんやら) 編

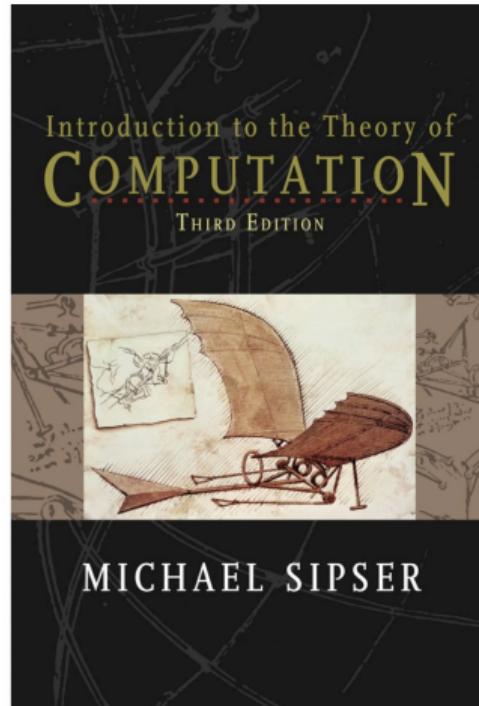
Hopcroft & Ullman,

*Introduction to automata theory
languages, and computation*



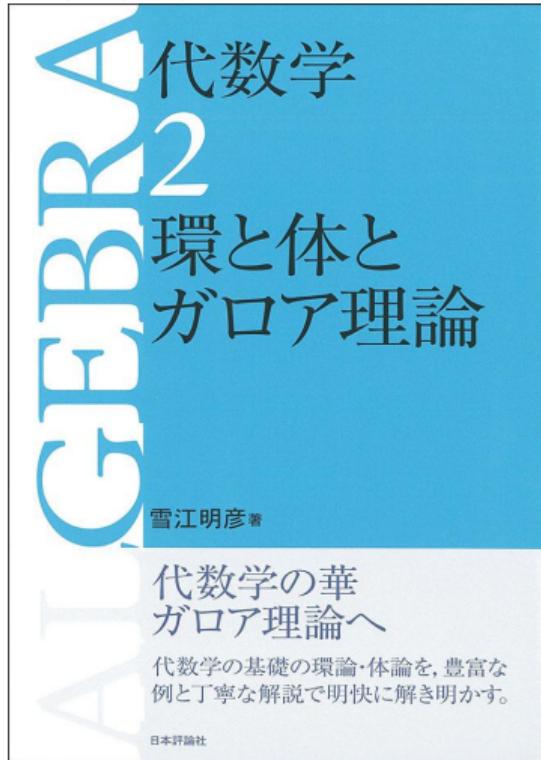
Sipser,

Introduction to the Theory of Computation



オススメの教科書: 有限体（ガロア体）と符号理論

工学的に良いのは知らないのですが、代数一般の視座からなら



第3章 体論の基本

- 3.1 体の拡大
- 3.2 代数閉包の存在
- 3.3 分離拡大
- 3.4 正規拡大
- 3.5 有限体
- 3.6 無限体上の多項式
- 3.7 単拡大

第4章 ガロア理論

- 4.1 ガロア拡大とガロアの基本定理

注: ガロア理論と有限体（ガロア体）は別物です

レポート課題: 上里担当分

レポート課題というより、演習問題を大問で4つ出します。全部解く想定です。
各問でセクションを区切り、合計4つのセクションからなる
PDFファイルを提出してください。

問1は、問2の準備です。 $\mathbb{F}[11]$ を構成してもらいます。

問2では、 $\mathbb{F}[11]$ を要素とする Vandermonde 行列を作り、
本当に逆行列が存在するかを手計算で確認してもらいます。
これで、簡易的な erasure coding ができたことになります。

問3は、演算数最適化を行なってください。

問4は、キャッシュ最適化を抽象化した pebble game に挑戦してください。

課題 1: $\mathbb{F}[11]$ の構成

要素数が 11 の有限体 $\mathbb{F}[11] = \{0, 1, \dots, 9, 10\}$ を構成したい。

足し算 $x \oplus y$ は普通の足し算をして 11 で割る: $(x + y) \bmod 11$

掛け算 $x \otimes y$ は普通の掛け算をして 11 で割る: $(x \times y) \bmod 11$

(1) まず足し算の逆元が必ず存在することを確認せよ
全ての x に対して、 $-x$ に相当するものを書け

(2) 同様に、掛け算の逆元が存在することを確認せよ
全ての非零 x に対して、 $\frac{1}{x}$ に相当する要素を全て書け

(3) 積の生成元 α で **最大のもの** を求めよ。積の生成元とは

$$\alpha, \alpha^2, \dots, \alpha^9, \alpha^{10}$$

が全て異なるもの。

α を掛けていくだけで、0 以外の全ての数を得られるので生成元と呼ぶ。

課題 2: 一般化 Vandermonde Matrix

高さ m 横幅 n 一般化 Vandermonde Matrix $V_{m,n}$ は 次で作れる:

$$V_{m,n} = m \begin{pmatrix} 1 & \alpha & \alpha^2 & \cdots & \alpha^{n-1} \\ 1 & (\alpha^2) & (\alpha^2)^2 & \cdots & (\alpha^2)^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & (\alpha^m) & (\alpha^m)^2 & \cdots & (\alpha^m)^{n-1} \end{pmatrix}^n$$

ただし、 $\alpha, \alpha^2, \alpha^3, \dots, \alpha^m$ は全て異なるとする。

(なので α としては積の生成元をとっておくと安心)

(1) $\mathbb{F}[11]$ を使う。課題 1 の (3) で求めた α を用いて

高さ (m) 4, 横幅 (n) 2 の Vandermonde matrix $V_{4,2}$ を生成し、それを記せ

(2) 好きな 2 行を削って 2×2 の正方行列を作り、それが逆行列を持つことを確かめよ

(これは課題ではないが、どの 2 行の組み合わせ (${}^4C_2 = 6$ 通り) で 2 行を削除をしても得られる 2×2 正方行列が正則となる。プログラムに計算させて確かめよ。また $V_{3,5}$ の行列を作り、どの 2 行削って得られる 3×3 行列も正則であることを確かめよ。)

課題 3: 速度最適化 1

以下のプログラムを \oplus の個数が最小になるように最適化せよ。
変数は幾つ追加しても良い。変数の個数は重要ではない。
とにかく字面に現れる \oplus の個数を減らせば良い。

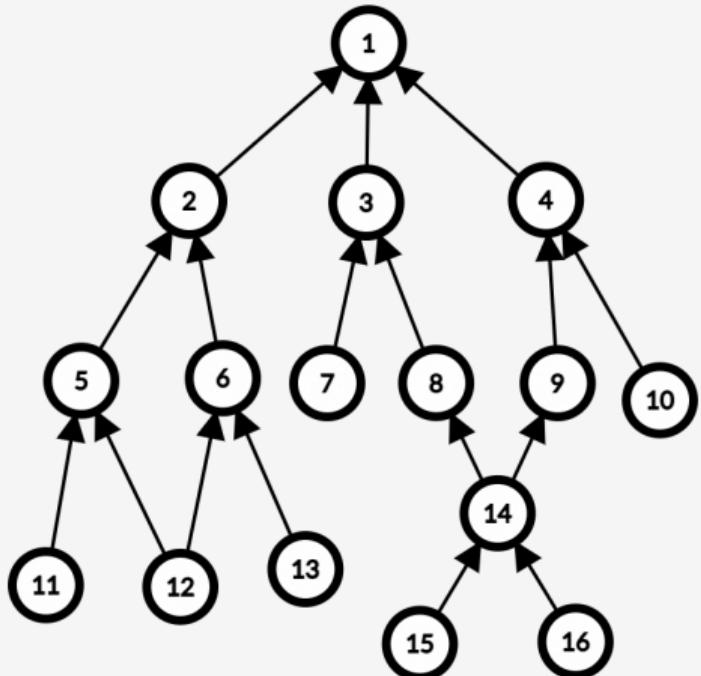
```
x1 ← a ⊕ b ⊕ c ⊕ w ⊕ x;  
x2 ← a ⊕ z ⊕ b ⊕ c ⊕ y;  
x3 ← c ⊕ a ⊕ b ⊕ y ⊕ w;  
x4 ← b ⊕ a ⊕ c ⊕ x ⊕ z;  
x5 ← b ⊕ c ⊕ a ⊕ w ⊕ z;  
return(x1, x2, x3, x4, x5);
```

ヒント 1: 可換則 $x \oplus y = y \oplus x$, 結合則 $x \oplus (y \oplus z) = (x \oplus y) \oplus z$ は勿論使うが,
XOR のキャンセル性質 $x \oplus x \oplus y = y$ は使わなくて良い。

ヒント 2: まず変数を一つ追加して $• \leftarrow • \oplus • \oplus •$ だけの形にしよう。

課題4: 速度最適化 2 Pebble Game

なるべく少ない pebble で解く = ①に pebble を置く プレイ列を（良い感じに）記せ。



(課題のための簡易版) ゲームのルール

- * 最初は全ての pebble が手元にある
- * 葉（入力辺がないノード）には
いつでも pebble が置ける
- * 内点（入力辺があるノード）には
その全ての直接の子が pebble を
持っている時に限り
 - 手元から置ける；または
 - どれかの子の pebble を移動させて良い
- * 1つのノードにおける pebble は一つまで
- * Pebble の回収には制約はない