

Distributed Storage Systems

Reliable Storage Continued:
Finite Fields & Linear Coding

Agenda



Reliable storage

Today's topics

- Basics of finite fields used in RAID and other storage systems
- Basics of coding for storage

Goals



This week's Learning Goals

- Understand finite field arithmetics
- Understand basics of coding for reliable storage

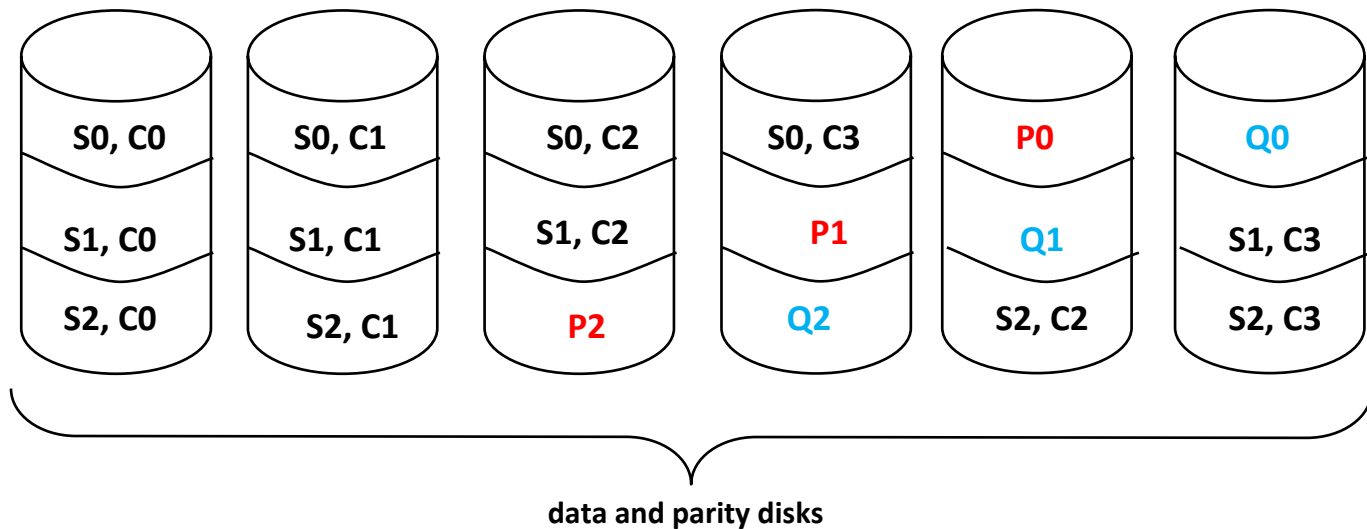
Class Structure

	Lecture	Lab
Week 1	Course introduction, networking basics, socket programming	Python sockets
Week 2	RPC, NFS, Practical RPC	Flask, JsonRPC, REST API
Week 3	AFS, reliable storage introduction	ZeroMQ, ProtoBuf
Week 4	Hard drives, RAID levels	RPi stack intro, RPi RAID with ZMQ
Week 5	Finite fields, Reed-Solomon Codes	Kodo intro, RS and RLNC with Kodo
Week 6	Repair problem, RS vs Regenerating codes	RPi simple distributed storage with Kodo RS
Week 7	Regenerating codes, XORBAS	RPi Regenerate lost fragments with RS
Week 8	Hadoop	RPi RLNC, recovery with recode
Week 9	Storage Virtualization, Network Attached Storage, Storage Area Networks	RPi basic HDFS (namenode+datanode, read and write pipeline)
Week 10	Object Storage	RPi basic S3 API
Week 11	Compression, Delta Encoding	Mini project consultation
Week 12	Data Deduplication	RPi Dedup
Week 13	Fog storage	Mini project consultation
Week 14	Security for Storage Systems and Recap	Mini project consultation

Reliable Storage

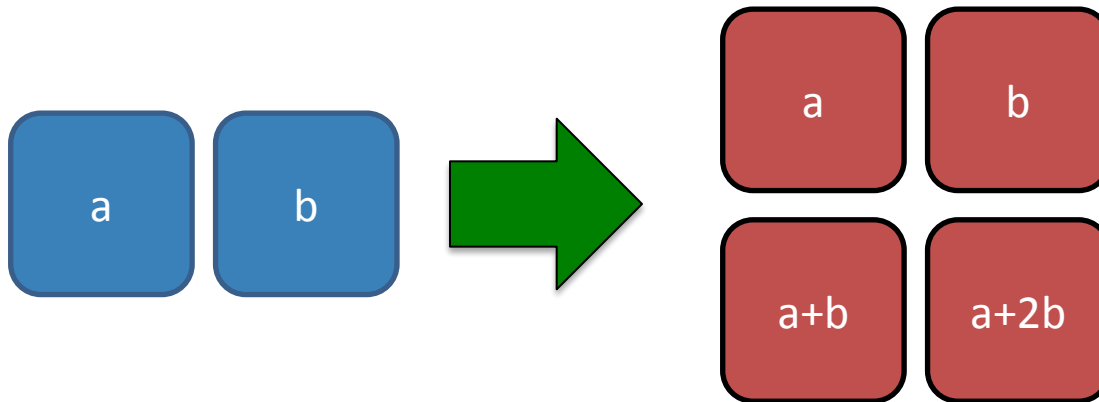
RAID6

- Level 5 with an extra parity
- Can tolerate two failures
- What are the odds of having two concurrent failures?



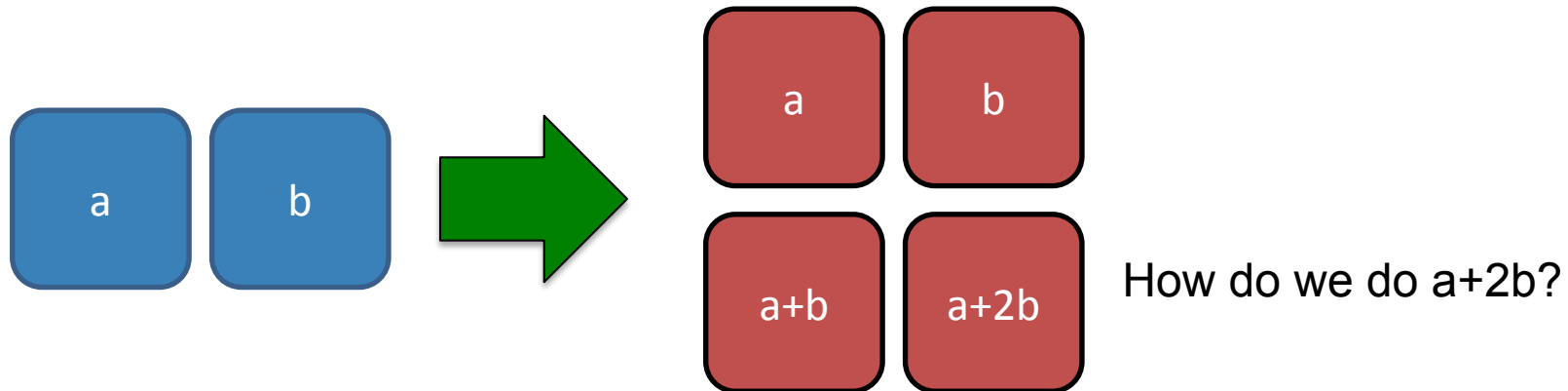
RAID6

- Level 5 with an extra parity
- Can tolerate two failures
- What are the odds of having two concurrent failures?



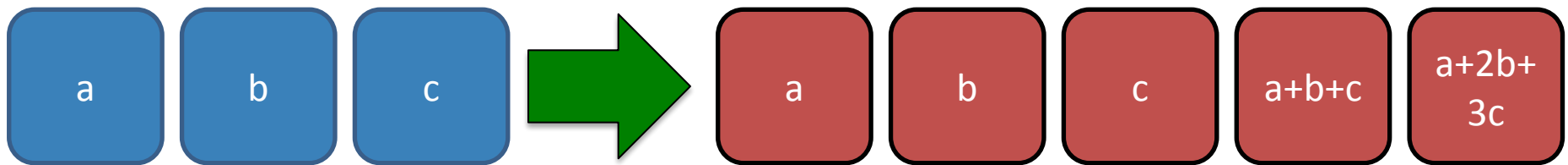
RAID6

- Level 5 with an extra parity
- Can tolerate two failures
- What are the odds of having two concurrent failures?



RAID6 - More than 2 stripes?

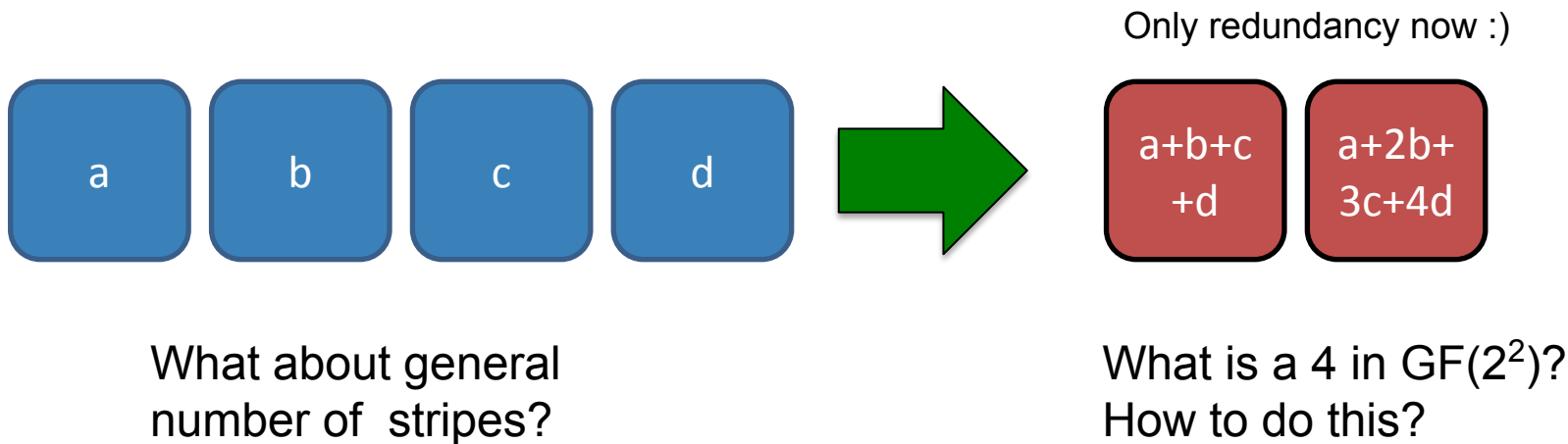
- Level 5 with an extra parity
- Can tolerate two failures



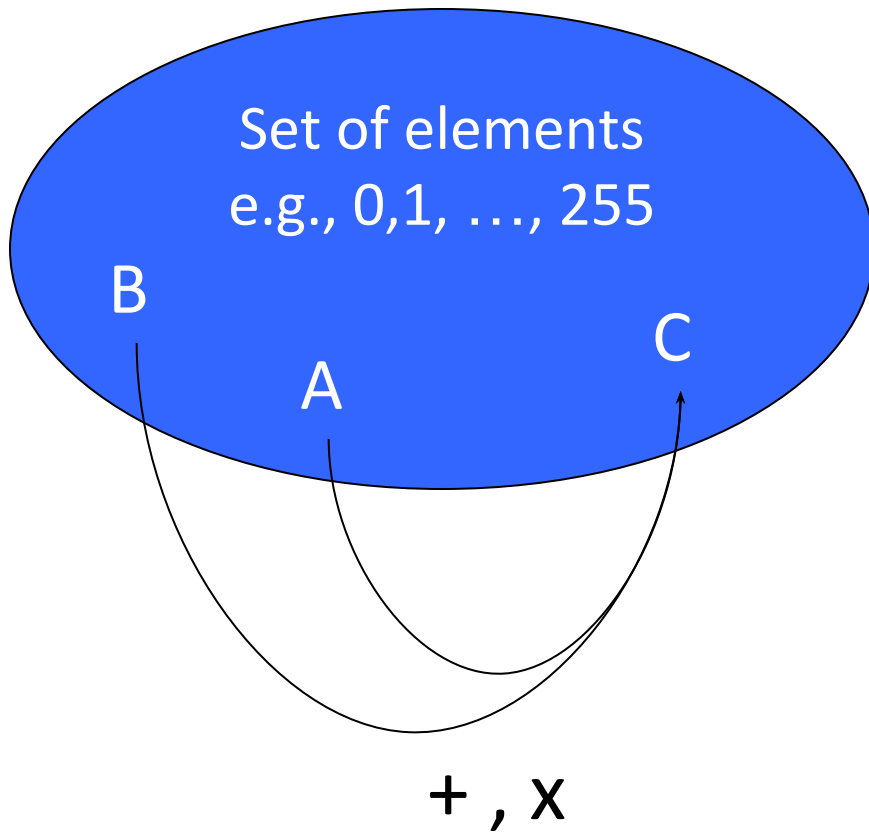
Reed Solomon codes provide the appropriate code construction

RAID6 - More than 2 stripes?

- Level 5 with an extra parity
- Can tolerate two failures



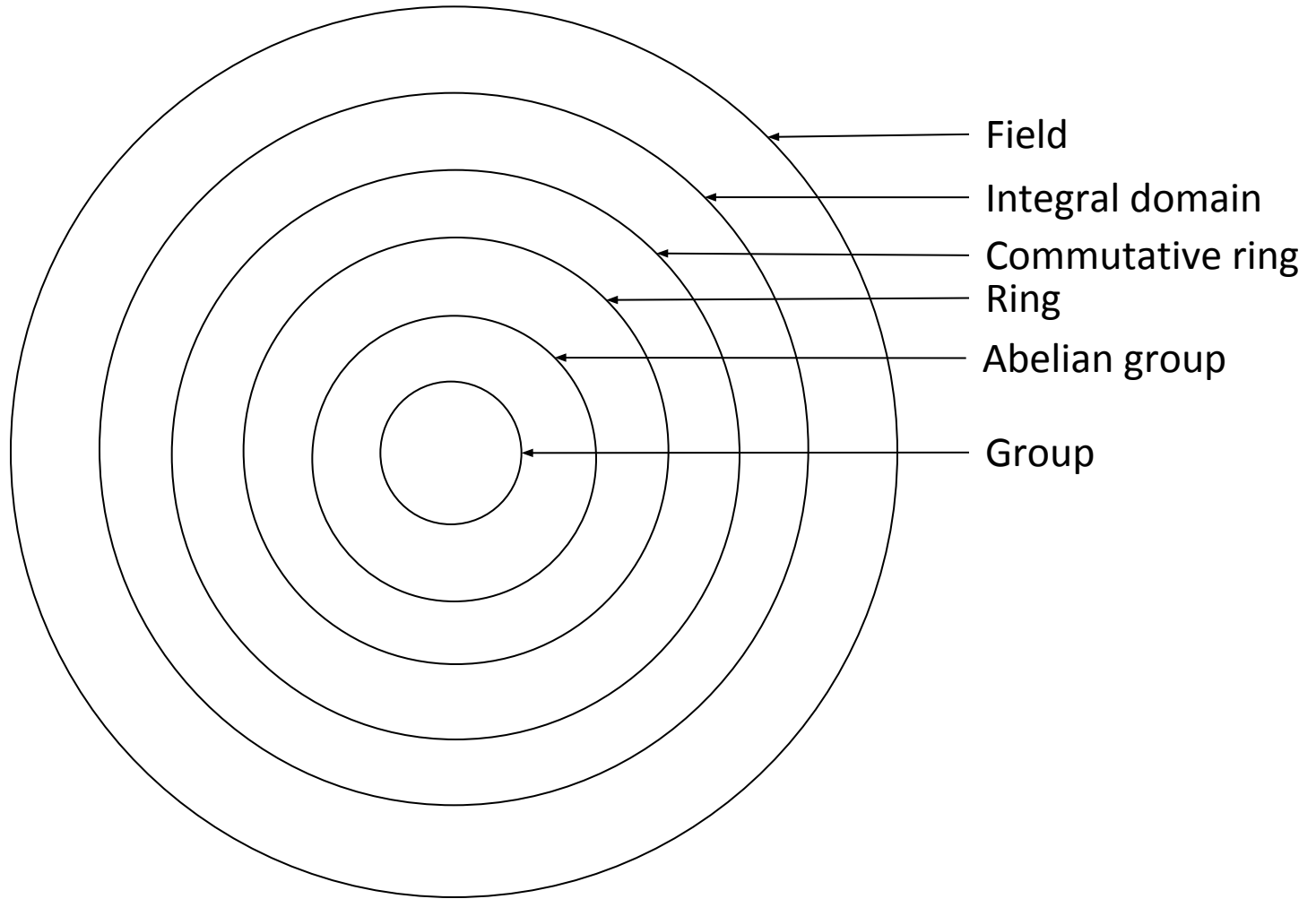
Finite Fields



Operations:
Addition, Multiplication

Property: closure

Groups, Rings, Fields



Groups, Rings, Fields

Groups

A group G , denoted $\{G, * \}$ is a set of elements with a binary operation $*$ that associates each ordered pair (a,b) of elements in G to an element $(a*b)$ in G following

Axioms

- *Closure*: If a and b in G , then $a * b$ is also in G
- *Associative*: $a * (b * c) = (a * b) * c$ for all a,b,c in G
- *Identity*: Exists e in G , s.t. $a*e = e*a = a$ for all a in G
- *Inverse*: for each a in G , exists a' in G , s.t.
 $a*a' = a'*a = e$

Finite group: finite number of elements

Groups, Rings, Fields

Abelian Group

Group that satisfies

- *Commutative*: If a and b in G , then $a * b = b * a$

Rings

A ring R , denoted by $\{R, +, \cdot\}$ is a set of elements with two binary operations: addition, multiplication. For all a, b, c in R the following axioms are satisfied

- R is **abelian group** with respect to the addition
- *Closure under multiplication*: ab in R
- *Associativity of multiplication*: $a(bc) = (ab)c$
- *Distributive laws*: $a(b+c) = ab + ac$

$$(a + b)c = ac + bc$$

Groups, Rings, Fields

Commutative Ring

A **ring** that also satisfies

- *Commutativity of multiplication*: $ab = ba$ in R

Integral Domain

R is a commutative **ring** that also satisfies

- *Multiplicative identity*: exists **1**, s.t. $a\mathbf{1} = \mathbf{1}a = a$
- *No zero divisors*: $ab = 0$, implies either $a = 0$ or $b = 0$

Field

F is a field, $\{F, +, \times\}$ that satisfies

- F is an integral domain
- *Multiplicative inverse*: for each a in F , except 0 , exists an element a^{-1} , s.t. $a(a^{-1}) = (a^{-1})a = 1$

Finite Fields GF(p)

Can write fields of the form $\text{GF}(p^n)$, where p is prime

Addition and multiplication over $\text{GF}(p)$ are mod p

Focus on $p = 2$

Example:

$\text{GF}(2)$ addition: equivalent to XOR

multiplication: equivalent to AND

How to divide? Multiply by multiplicative inverse

Finding the multiplicative inverse

1.- Can look for a^{-1} such that $(a^{-1} \cdot a) \equiv 1$

2.- Can use the extended Euclidean algorithm

Finite Fields - Applying GF(2) to NC

Example:

GF(2) addition: XOR

multiplication: AND

Given 2 data packets

P1: 01011001 P2: 10001001

calculate the content of the coded packet $P1+P2$.

What are the coefficients?

$$\begin{array}{rcl} P1 + P2 = & 01011001 & \text{(XOR bit by bit)} \\ & 10001001 & \\ & 11010000 & \end{array}$$

Modular Arithmetic

Modulus

If a is an integer, $n > 0$ integer, we define $a \bmod n$ to be the remainder when a is divided by n

- The integer n is called the *modulus*
- For any integer a , we can write

$$a = qn + r, \text{ with } 0 \leq r < n, \text{ and } q = \lfloor a/n \rfloor$$

- E.g., $11 \bmod 7 = 4$, $-11 \bmod 7 = 3$

Congruent modulo n

If $(a \bmod n) = (b \bmod n)$, and it's expressed

$$a \equiv b \pmod{n}$$

- E.g., $20 \equiv 6 \pmod{7}$

Modular Arithmetic

Properties of congruencies

- $a \equiv b \pmod{n}$ if $n \mid (a-b)$
- $a \equiv b \pmod{n}$ implies $b \equiv a \pmod{n}$
- $a \equiv b \pmod{n}$ and $b \equiv c \pmod{n}$,
then $a \equiv c \pmod{n}$

Modular arithmetic operations

- $[(a \bmod n) + (b \bmod n)] \bmod n = (a+b) \bmod n$
- $[(a \bmod n) - (b \bmod n)] \bmod n = (a-b) \bmod n$
- $[(a \bmod n) \times (b \bmod n)] \bmod n = (a \times b) \bmod n$

Rules of ordinary arithmetic involving addition,
subtraction, multiplication carry over

Modular Arithmetic

Properties of modular arithmetic

Define $Z_n = \{ 0, 1, \dots, n-1 \}$ as the set of **residues** or **residue classes** mod n .

Each element of Z_n is a residue class and can define it as

$$[j] = \{a : a \text{ is integer, } a \equiv j \pmod{n}\}$$

- Reducing k mod n : finding smallest non-negative integer a , such that $k \equiv a \pmod{n}$

Modular Arithmetic

Properties of modular arithmetic

- $(w + x) \bmod n = (x + w) \bmod n$
- $(w \times y) \bmod n = (y \times w) \bmod n$
- $((w + x) + y) \bmod n = (w + (x + y)) \bmod n$
- $((w \times j) \times y) \bmod n = (w \times (j \times y)) \bmod n$
- $(w \times (y + j)) \bmod n = ((w \times y) + (w \times j)) \bmod n$
- $(0 + w) \bmod n = w \bmod n$
- $(1 \times w) \bmod n = w \bmod n$

Modular Arithmetic

Properties of modular arithmetic

- If $(a + b) \equiv (a + c) \pmod{n}$, then $b \equiv c \pmod{n}$
- If $(a \times b) \equiv (a \times c) \pmod{n}$, then $b \equiv c \pmod{n}$ if a is relatively prime to n , i.e., $\gcd(a, n) = 1$

Why is the last property important?

What happens when “p” is not a prime?

Will modular arithmetic still work?

Example 1:

If $\gcd(a,n) \neq 1$, the last equation does not hold

e.g. $6 \times 3 = 18 = 2 \pmod{8}$

and $6 \times 7 = 42 = 2 \pmod{8}$ but

$3 \pmod{8} \neq 7 \pmod{8}$

What happens when “p” is not a prime?

Will modular arithmetic still work?

+	0	1	2	3
0	0	1	2	3
1	1	2	3	0
2	2	3	0	1
3	3	0	1	2

x	0	1	2	3
0	0	0	0	0
1	0	1	2	3
2	0	2	0	2
3	0	3	2	1

What about GF(2ⁿ)?

Since 2ⁿ is not a prime, operations are defined in a different way => **polynomial arithmetic**

Ordinary polynomial arithmetic:

A polynomial of degree n

$$F(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0 x^0 = \sum a_i x^i$$

a_i are the coefficients, chosen from a set

Operations:

Addition $f(x) + g(x) = \sum (a_i + b_i) x^i$

Multiplication $f(x) \times g(x) = \sum C_i x^i$

with $c_k = a_0 b_k + a_1 b_{k-1} + \dots + a_k b_0$

What about $\text{GF}(2^n)$?

Polynomial arithmetic in $\text{GF}(2^n)$:

- Arithmetic follows rules of polynomial arithmetic
- Arithmetic of coefficients is performed modulo 2
 - i.e., using $\text{GF}(2)$ addition/multiplication for coefficients of the same order
 - e.g., $(a_i + b_i) \bmod 2$
- If multiplication results in a polynomial greater than $n-1$, then the polynomial is reduced modulo an irreducible polynomial $p(x)$
 - Think of it as a $\bmod p(x)$ operation: divide by $p(x)$, keep the remainder

Example GF(2²)

Irreducible polynomial $p(x) = x^2 + x + 1 \quad (111)_b$

+	0	1	2	3
0	0	1	2	3
1	1	0	3	2
2	2	3	0	1
3	3	2	1	0

How about $2 + 3$?

$2 = (10)_b$ and $3 = (11)_b$

As polynomials:

$2 \equiv x$ and $3 \equiv x + 1$

Thus, $2 + 3$ becomes

$x + (x + 1) = 1$

Example GF(2²)

Irreducible polynomial $p(x) = x^2 + x + 1 \quad (111)_b$

Equal number of each element: RLNC's properties
are maintained when recoding

How about 2×3 ?

$$x(x+1) = (x^2 + x) \bmod p(x)$$

Multiplicative inverses; easy to spot in table

Can we compute without generating table?

$$\begin{array}{r} x^2 + x \\ x^2 + x + 1 \\ \hline 1 \end{array}$$

How to implement multiplication $GF(2^n)$?

A.- Product (shifts+ XORs)

- 1) Pick one as multiplier (M) and another as multiplicand (m)
- 2) For each “1” in “M”, left shift the “m” by the position of the “1”
- 3) XOR shifted versions

B.- Modulo irreducible polynomial (long division)

- 4) Initialize: $F(0)$ = Result of part A
- 5) Take irreducible polynomial $p(x)$ left shift until first “1” of polynomial and of value F match
- 6) $F(i+1) = (\text{shifted } p(x)) \text{ XOR } (F(i))$
- 7) Stop if first “1” of $F(i+1)$ occurs in the $(n-1)$ -th bit

How to implement multiplication?

$$M = (00010001)_b \text{ and } m = (10100111)_b$$

$$\begin{array}{r} 101001110000 \quad (m \text{ shifted 4 times}) \\ (XOR) \quad 10100111 \quad (m \text{ shifted 0 times}) \\ \hline 101011010111 \end{array}$$

B.- Modulo irreducible polynomial

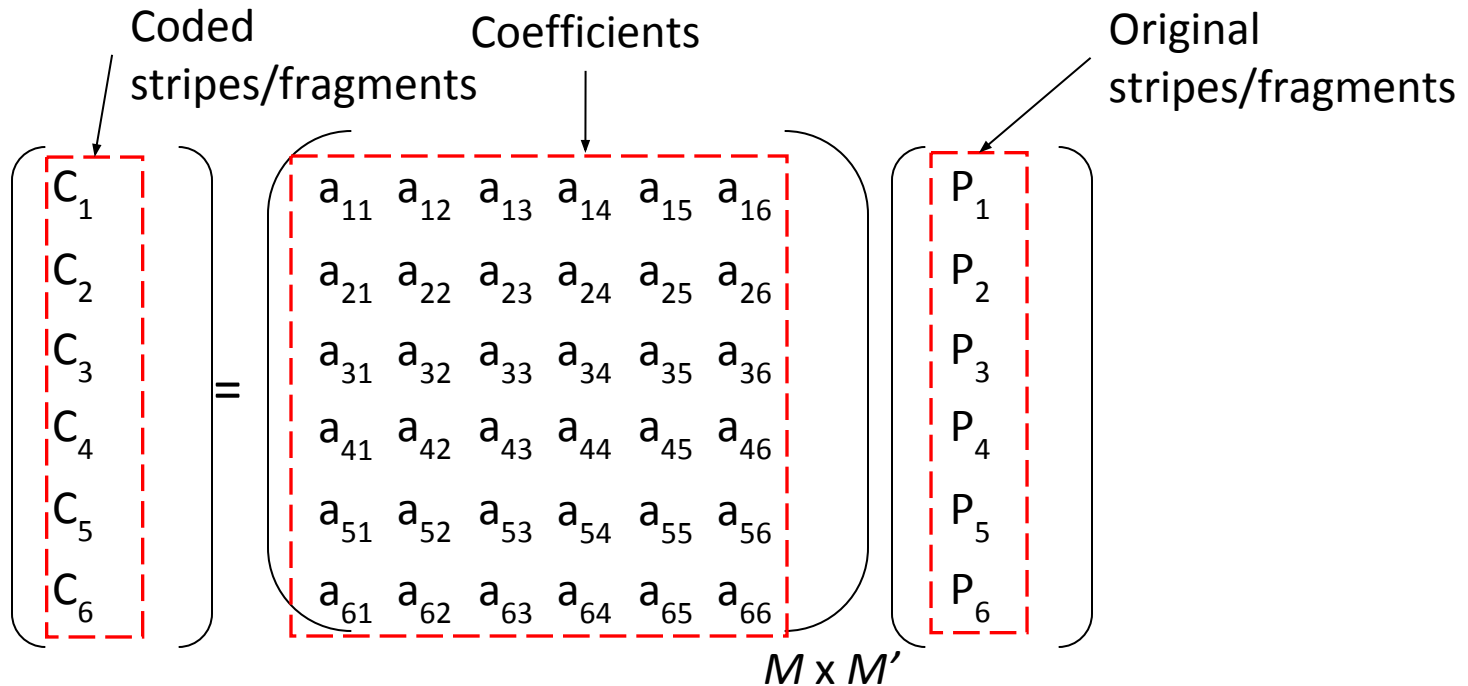
$$p(x) = x^8 + x^4 + x^3 + x + 1 \quad (100011011)_b$$

$$\begin{array}{r} 101011010111 \quad \text{mod } p(x) \\ (XOR) \quad \mathbf{100011011000} \\ \hline 001000001111 \\ (XOR) \quad 00\mathbf{1000110110} \\ \hline 000000\mathbf{111001} \rightarrow (\mathbf{00111001})_b \end{array}$$

Generating Coded Fragments/Stripes

- Generating a linear coded fragment/stripe (C)

$$C_i = \sum a_{ij} P_j$$



How you pick the coefficients determines properties and performance

Generating Coded Packets

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} & a_{15} & a_{16} \\ a_{21} & a_{22} & a_{23} & a_{24} & a_{25} & a_{26} \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} & a_{36} \\ a_{41} & a_{42} & a_{43} & a_{44} & a_{45} & a_{46} \\ a_{51} & a_{52} & a_{53} & a_{54} & a_{55} & a_{56} \\ a_{61} & a_{62} & a_{63} & a_{64} & a_{65} & a_{66} \end{pmatrix} \begin{pmatrix} p_{11} & p_{12} & p_{13} & p_{14} & p_{15} \\ p_{21} & p_{22} & p_{23} & p_{24} & p_{25} \\ p_{31} & p_{32} & p_{33} & p_{34} & p_{35} \\ p_{41} & p_{42} & p_{43} & p_{44} & p_{45} \\ p_{51} & p_{52} & p_{53} & p_{54} & p_{55} \\ p_{61} & p_{62} & p_{63} & p_{64} & p_{65} \end{pmatrix}$$

RAID6

- Some data coded, some is not
- Called: 'Systematic' Code

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ a_{11} & a_{12} & a_{13} & a_{14} & a_{15} & a_{16} \\ a_{21} & a_{22} & a_{23} & a_{24} & a_{25} & a_{26} \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$

RAID6

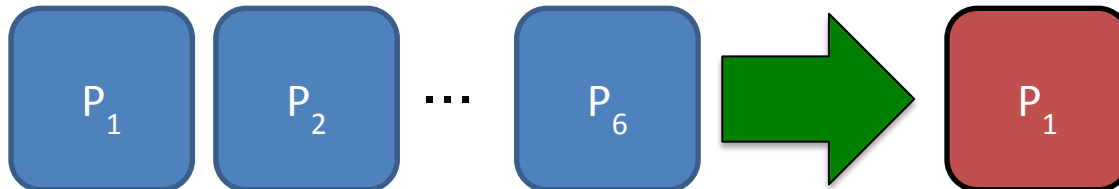
- Some data coded, some is not
- Called: 'Systematic' Code

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$

RAID6

- Some data coded, some is not
- Called: 'Systematic' Code

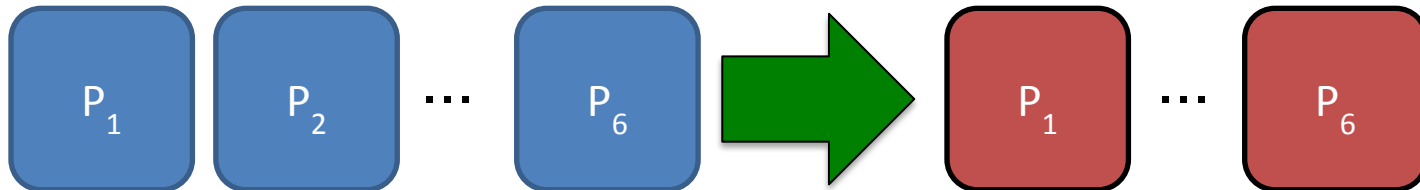
$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{bmatrix}$$



RAID6

- Some data coded, some is not
- Called: 'Systematic' Code

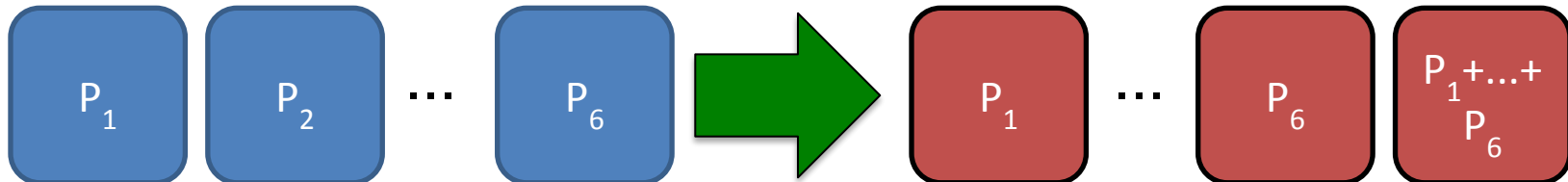
$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{bmatrix}$$



RAID6

- Some data coded, some is not
- Called: 'Systematic' Code

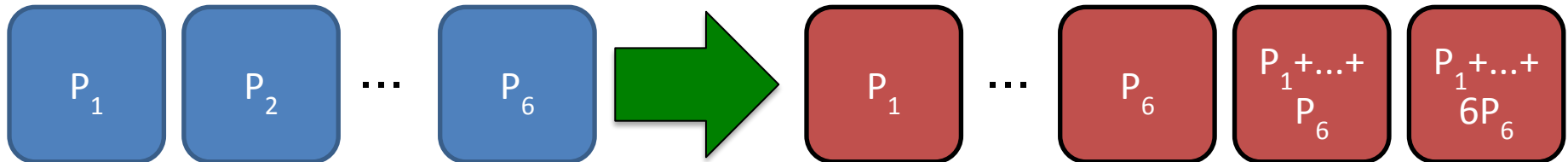
$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{bmatrix}$$



RAID6

- Some data coded, some is not
- Called: 'Systematic' Code

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{bmatrix}$$



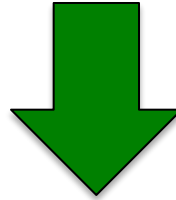
RAID6 - Decode

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{bmatrix}$$

The matrix is an 8x6 matrix. The first four rows are the identity matrix. The fifth and sixth rows are the parity rows, with the last two columns being 1 and 0, and 0 and 1 respectively. The seventh and eighth rows are the data rows, with the last two columns being 1 and 1, and 2 and 3 respectively. The matrix is partitioned into two 4x3 blocks by a dashed red line between the fourth and fifth rows. The first block contains the identity matrix and the parity rows. The second block contains the data rows and the parity rows.

RAID6 - Decode

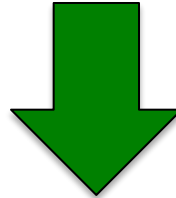
$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix}$$

RAID6 - Decode

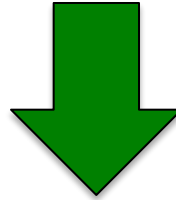
$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 3 & 4 & 5 & 6 \end{pmatrix}$$

RAID6 - Decode

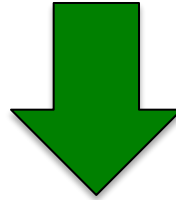
$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 5 & 6 \end{pmatrix}$$

RAID6 - Decode

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$

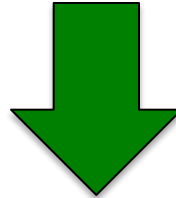


**Multiply by 5
and
Add**

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 5 & 6 \end{pmatrix}$$

RAID6 - Decode

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$

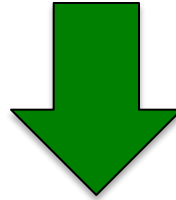



$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{pmatrix}$$

Multiply by 3^{-1}

RAID6 - Decode

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$

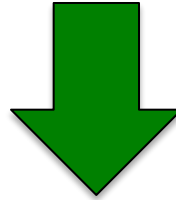


Add 

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

RAID6 - Decode

$$\begin{pmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_7 \\ C_8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix} \begin{pmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

Decoded

**What are good
constructions of the
(coefficient) matrix?**

What are good constructions of the (coefficient) matrix?

Depends on performance objective:

- Storage reduction
- Reduced I/O operations
- ...

What are good constructions of the (coefficient) matrix?

Depends on performance objective:

- **Storage reduction:** MDS codes, e.g., Reed Solomon
- Reduced I/O operations
- ...

MDS Codes

- A **maximum distance separable** code, denoted by **MDS(n, k)**, has the property that any **k** ($< n$) out of n nodes can be used to reconstruct original native blocks
 - i.e., at most $n-k$ disk failures can be tolerated
- Example:
 - RAID-5 is an MDS($n, n-1$) code as it can tolerate 1 disk failure
 - RAID-6 is an MDS($n, n-2$) code as it can tolerate 2 disk failures

Reed Solomon - Vandermonde

- Not systematic

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ \vdots \\ C_n \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ \alpha_1^1 & \alpha_2^1 & \alpha_3^1 & \alpha_4^1 & \dots & \alpha_k^1 \\ \alpha_1^2 & \alpha_2^2 & \alpha_3^2 & \alpha_4^2 & \dots & \alpha_k^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \alpha_1^{(n-1)} & \alpha_2^{(n-1)} & \alpha_3^{(n-1)} & \alpha_4^{(n-1)} & \dots & \alpha_k^{(n-1)} \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ \vdots \\ P_k \end{bmatrix}$$

Reed Solomon - Vandermonde

- Not systematic

$$H = \begin{pmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ \alpha_1^1 & \alpha_2^1 & \alpha_3^1 & \alpha_4^1 & \dots & \alpha_k^1 \\ \alpha_1^2 & \alpha_2^2 & \alpha_3^2 & \alpha_4^2 & \dots & \alpha_k^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \alpha_1^{(n-1)} & \alpha_2^{(n-1)} & \alpha_3^{(n-1)} & \alpha_4^{(n-1)} & \dots & \alpha_k^{(n-1)} \end{pmatrix}$$

Reed Solomon - Vandermonde

- Not systematic

$$H = \begin{pmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & 4 & \dots & k \\ 1^2 & 2^2 & 3^2 & 4^2 & \dots & k^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1^{(n-1)} & 2^{(n-1)} & 3^{(n-1)} & 4^{(n-1)} & \dots & k^{(n-1)} \end{pmatrix}$$

There is a limit on the value of $n \rightarrow n$ cannot be larger than (size of field -1)
For $GF(2^m)$, $n < 2^m - 1$

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$

$$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 2^2 & 3^2 & 4^2 \\ 1 & 2^3 & 3^3 & 4^3 \\ 1 & 2^4 & 3^4 & 4^4 \\ 1 & 2^5 & 3^5 & 4^5 \end{pmatrix}$$

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - polynomial: x^3+x^2+1

$$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 2^2 & 3^2 & 4^2 \\ 1 & 2^3 & 3^3 & 4^3 \\ 1 & 2^4 & 3^4 & 4^4 \\ 1 & 2^5 & 3^5 & 4^5 \end{pmatrix} \quad \begin{array}{l} \text{Mapping to polynomial} \\ 0 \rightarrow 0 \\ 1 \rightarrow 1 \\ 2 \rightarrow x \\ 3 \rightarrow x+1 \\ 4 \rightarrow x^2 \\ 5 \rightarrow x^2+1 \\ 6 \rightarrow x^2+x \\ 7 \rightarrow x^2+x+1 \end{array}$$

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 2^2 & 3^2 & 4^2 \\ 1 & 2^3 & 3^3 & 4^3 \\ 1 & 2^4 & 3^4 & 4^4 \\ 1 & 2^5 & 3^5 & 4^5 \end{pmatrix}$	<p>Mapping to polynomial</p> <p> $0 \rightarrow 0$ $1 \rightarrow 1$ $2 \rightarrow x$ $3 \rightarrow x + 1$ $4 \rightarrow x^2$ $5 \rightarrow x^2 + 1$ $6 \rightarrow x^2 + x$ $7 \rightarrow x^2 + x + 1$ </p>	<p>Mapping to polynomial</p> <p> $2^2 \rightarrow x^2 \rightarrow 4$ $2^3 \rightarrow x^3 \rightarrow R(x^3/x^3+x^2+1) \rightarrow x^2+1 \rightarrow 5$ $2^4 \rightarrow x^4 \rightarrow x(x^2+1) \rightarrow x^2+x+1 \rightarrow 7$ $2^5 \rightarrow x.x^4 \rightarrow x(x^2+x+1) \rightarrow x+1 \rightarrow 3$ $2^6 \rightarrow x.x^5 \rightarrow x(x+1) \rightarrow x^2+x \rightarrow 6$ $2^7 \rightarrow x.x^6 \rightarrow x(x^2+x) \rightarrow x^3+x^2 \rightarrow 1$ </p>
--	--	---

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 3^2 & 4^2 \\ 1 & 5 & 3^2 & 4^2 \\ 1 & 7 & 3^2 & 4^2 \\ 1 & 3 & 3^5 & 4^5 \end{pmatrix}$	<p>Mapping to polynomial</p> <p> $0 \rightarrow 0$ $1 \rightarrow 1$ $2 \rightarrow x$ $3 \rightarrow x + 1$ $4 \rightarrow x^2$ $5 \rightarrow x^2 + 1$ $6 \rightarrow x^2 + x$ $7 \rightarrow x^2 + x + 1$ </p>	<p>Mapping to polynomial</p> <p> $2^2 \rightarrow x^2 \rightarrow \mathbf{4}$ $2^3 \rightarrow x^3 \rightarrow R(x^3/x^3+x^2+1) \rightarrow x^2+1 \rightarrow \mathbf{5}$ $2^4 \rightarrow x^4 \rightarrow x(x^2+1) \rightarrow x^2+x+1 \rightarrow \mathbf{7}$ $2^5 \rightarrow x.x^4 \rightarrow x(x^2+x+1) \rightarrow x+1 \rightarrow \mathbf{3}$ $2^6 \rightarrow x.x^5 \rightarrow x(x+1) \rightarrow x^2+x \rightarrow \mathbf{6}$ $2^7 \rightarrow x.x^6 \rightarrow x(x^2+x) \rightarrow x^3+x^2 \rightarrow \mathbf{1}$ </p>
--	--	---

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 3^2 & 4^2 \\ 1 & 5 & 3^3 & 4^3 \\ 1 & 7 & 3^4 & 4^4 \\ 1 & 3 & 3^5 & 4^5 \end{pmatrix}$	<p>Mapping to polynomial</p> <p> $0 \rightarrow 0$ $1 \rightarrow 1$ $2 \rightarrow x$ $3 \rightarrow x + 1$ $4 \rightarrow x^2$ $5 \rightarrow x^2 + 1$ $6 \rightarrow x^2 + x$ $7 \rightarrow x^2 + x + 1$ </p>	<p>Mapping to polynomial</p> <p> $3^2 \rightarrow (x+1)^2 \rightarrow x^2 + 1 \rightarrow \mathbf{5}$ $3^3 \rightarrow (x^2+1)(x+1) \rightarrow x \rightarrow \mathbf{2}$ $3^4 \rightarrow x(x+1) \rightarrow (x^2+x) \rightarrow \mathbf{6}$ $3^5 \rightarrow (x^2+x)(x+1) \rightarrow x^2 + x + 1 \rightarrow \mathbf{7}$ </p>
--	--	---

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 4^2 \\ 1 & 5 & 2 & 4^3 \\ 1 & 7 & 6 & 4^4 \\ 1 & 3 & 7 & 4^5 \end{pmatrix}$$

Mapping to polynomial	Mapping to polynomial	
$0 \rightarrow 0$	$3^2 \rightarrow (x+1)^2 \rightarrow x^2+1$	$\rightarrow 5$
$1 \rightarrow 1$	$3^3 \rightarrow (x^2+1)(x+1) \rightarrow x$	$\rightarrow 2$
$2 \rightarrow x$	$3^4 \rightarrow x(x+1) \rightarrow (x^2+x)$	$\rightarrow 6$
$3 \rightarrow x+1$	$3^5 \rightarrow (x^2+x)(x+1) \rightarrow x^2+x+1$	$\rightarrow 7$
$4 \rightarrow x^2$		
$5 \rightarrow x^2+1$		
$6 \rightarrow x^2+x$		
$7 \rightarrow x^2+x+1$		

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 4^2 \\ 1 & 5 & 2 & 4^3 \\ 1 & 7 & 6 & 4^4 \\ 1 & 3 & 7 & 4^5 \end{pmatrix}$	<p>Mapping to polynomial</p> <p>$0 \rightarrow 0$</p> <p>$1 \rightarrow 1$</p> <p>$2 \rightarrow x$</p> <p>$3 \rightarrow x + 1$</p> <p>$4 \rightarrow x^2$</p> <p>$5 \rightarrow x^2 + 1$</p> <p>$6 \rightarrow x^2 + x$</p> <p>$7 \rightarrow x^2 + x + 1$</p>	<p>Mapping to polynomial</p> <p>$4^2 \rightarrow x^4 \rightarrow x^2 + x + 1 \rightarrow \mathbf{7}$</p> <p>$4^3 \rightarrow x^2(x^2 + x + 1) \rightarrow x^2 + x \rightarrow \mathbf{6}$</p> <p>$4^4 \rightarrow x^2(x^2 + x) \rightarrow x \rightarrow \mathbf{2}$</p> <p>$4^5 \rightarrow x^3 \rightarrow x^2 + 1 \rightarrow \mathbf{5}$</p>
--	--	--

Reed Solomon - Vandermonde

- Not systematic $n = 6, k = 4$
- Need $GF(2^3)$ (at least) - primitive polynomial: x^3+x^2+1
- What about the exponents?

$H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 7 \\ 1 & 5 & 2 & 6 \\ 1 & 7 & 6 & 2 \\ 1 & 3 & 7 & 5 \end{pmatrix}$	<p>Mapping to polynomial</p> <p>$0 \rightarrow 0$</p> <p>$1 \rightarrow 1$</p> <p>$2 \rightarrow x$</p> <p>$3 \rightarrow x + 1$</p> <p>$4 \rightarrow x^2$</p> <p>$5 \rightarrow x^2 + 1$</p> <p>$6 \rightarrow x^2 + x$</p> <p>$7 \rightarrow x^2 + x + 1$</p>	<p>Mapping to polynomial</p> <p>$4^2 \rightarrow x^4 \rightarrow x^2 + x + 1 \rightarrow \mathbf{7}$</p> <p>$4^3 \rightarrow x^2(x^2 + x + 1) \rightarrow x^2 + x \rightarrow \mathbf{6}$</p> <p>$4^4 \rightarrow x^2(x^2 + x) \rightarrow x \rightarrow \mathbf{2}$</p> <p>$4^5 \rightarrow x^3 \rightarrow x^2 + 1 \rightarrow \mathbf{5}$</p>
--	--	--

Reed Solomon - Vandermonde

- Systematic - different procedures starting with H
- Simple:
 - Transpose H
 - Perform Gaussian Elimination on $H^T \rightarrow H_{\text{sys}}^T = [\mathbf{I} \mid \mathcal{R}]$
 - Transpose again to reach H_{sys}

Reed Solomon - Cauchy

- Systematic by design

$$H = \begin{pmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \\ (x_1+y_1)^{-1} & (x_2+y_1)^{-1} & (x_3+y_1)^{-1} & (x_4+y_1)^{-1} & \dots & (x_k+y_1)^{-1} \\ (x_1+y_2)^{-1} & (x_2+y_2)^{-1} & (x_3+y_2)^{-1} & (x_4+y_2)^{-1} & \dots & (x_k+y_2)^{-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ (x_1+y_{n-k})^{-1} & (x_2+y_{n-k})^{-1} & (x_3+y_{n-k})^{-1} & (x_4+y_{n-k})^{-1} & \dots & (x_k+y_{n-k})^{-1} \end{pmatrix}$$

$X = \{x_1, x_2, x_3, \dots, x_k\}$ and $Y = \{y_1, y_2, y_3, \dots, y_{n-k}\}$, where $x_i \neq y_j \forall i, j$

Reed Solomon - Cauchy

- Systematic by design

$H =$

$$\begin{pmatrix}
 1 & 0 & 0 & 0 & \dots & 0 \\
 0 & 1 & 0 & 0 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & 0 & 0 & 0 & \dots & 1 \\
 (x_1+y_1)^{-1} & (x_2+y_1)^{-1} & (x_3+y_1)^{-1} & (x_4+y_1)^{-1} & \dots & (x_k+y_1)^{-1} \\
 (x_1+y_2)^{-1} & (x_2+y_2)^{-1} & (x_3+y_2)^{-1} & (x_4+y_2)^{-1} & \dots & (x_k+y_2)^{-1} \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 (x_1+y_{n-k})^{-1} & (x_2+y_{n-k})^{-1} & (x_3+y_{n-k})^{-1} & (x_4+y_{n-k})^{-1} & \dots & (x_k+y_{n-k})^{-1}
 \end{pmatrix}$$

This addition
and inverse is in
GF(q)



Example: $k = 2, n = 7 \rightarrow X = \{1, 2\}$ and $Y = \{0, 3, 4, 5, 6\}$

Reed Solomon - Cauchy

- $n = 6$, $k = 4$ and $X = \{1, 2, 3, 4\}$ and $Y = \{0, 5\}$

$$H = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 2^{-1} & 3^{-1} & 4^{-1} \\ 4^{-1} & 7^{-1} & 6^{-1} & 1 \end{pmatrix}$$

Find inverses
How?

Reed Solomon - Cauchy

- $n = 6$, $k = 4$ and $X = \{1, 2, 3, 4\}$ and $Y = \{0, 5\}$

$$H = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 2^{-1} & 3^{-1} & 4^{-1} \\ 4^{-1} & 7^{-1} & 6^{-1} & 1 \end{pmatrix}$$

Tables or Computation

$$2 \cdot 2^{-1} = 1 \quad ?$$

$$2 \cdot 6 = x(x^2 + x) \rightarrow 1$$

$$\text{Thus, } 2^{-1} = 6$$

Reed Solomon - Cauchy

- $n = 6$, $k = 4$ but with $\mathbf{X} = \{1, 3, 4, 5\}$ and $\mathbf{Y} = \{0, 2\}$

$$\mathbf{H} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 3^{-1} & 4^{-1} & 5^{-1} \\ 3^{-1} & 1 & 6^{-1} & 7^{-1} \end{pmatrix}$$

Large number of Cauchy
Alternatives - not all built equally

Distributed Storage: Reliability Beyond a Single Server

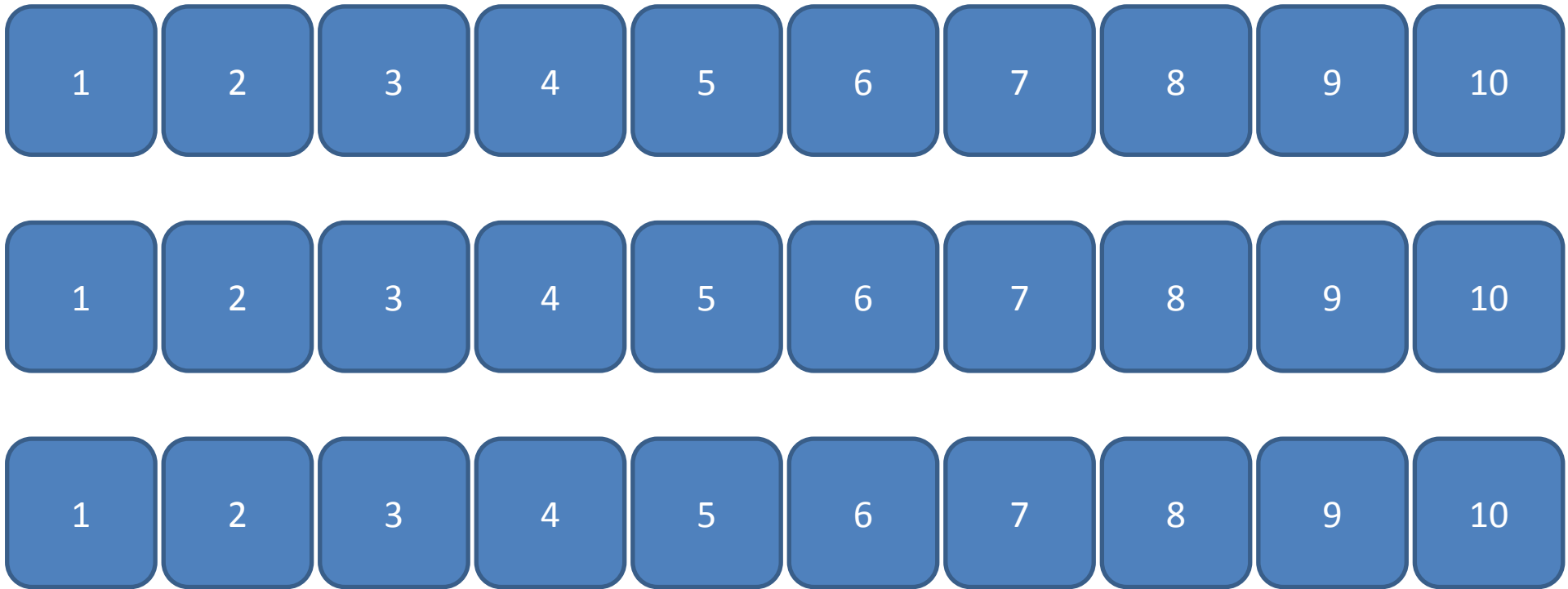


Distributed storage systems

- Numerous disk failures per day.
- Failures are the norm rather than the exception
- Must introduce **redundancy for reliability**
- Replication or erasure coding?
 - Current question in Hadoop HDFS, OpenStack Swift, Ceph, ...

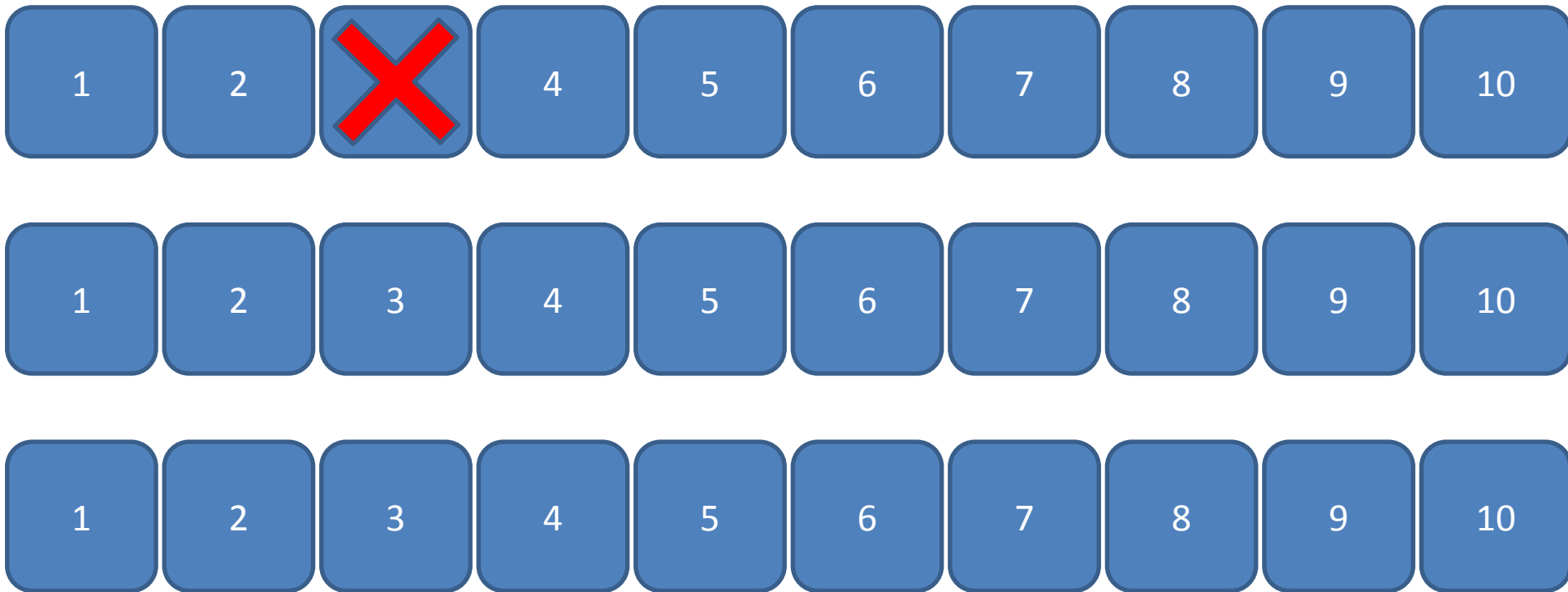


State of the Art



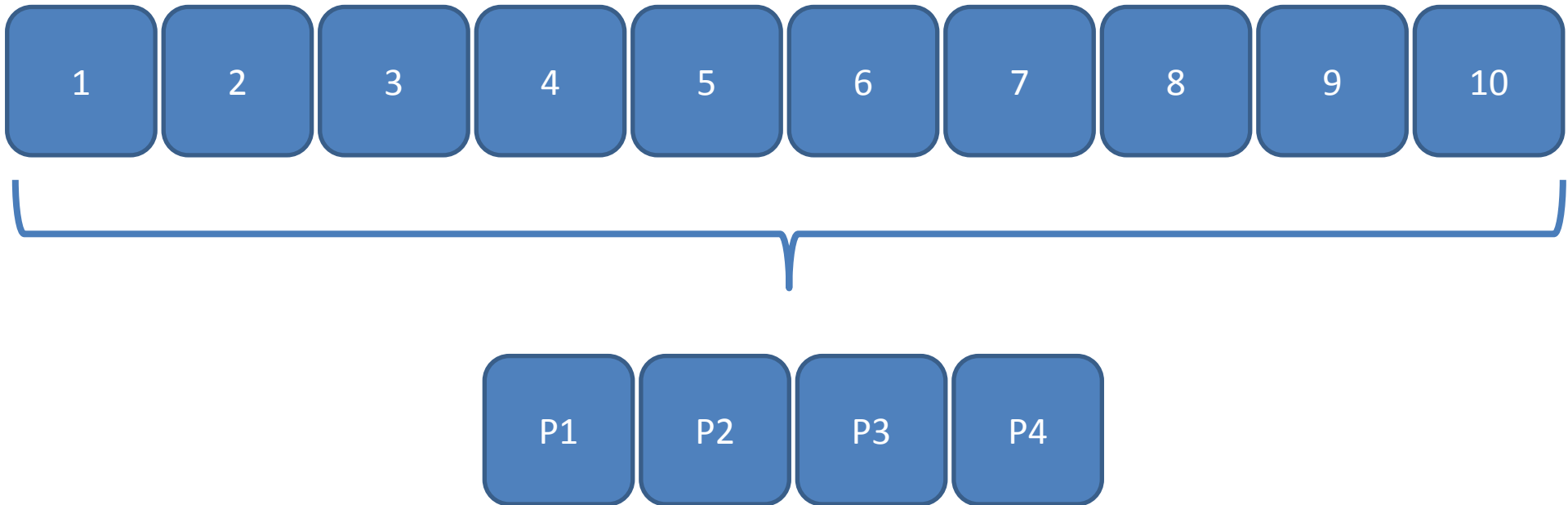
Overhead: 200%

State of the Art



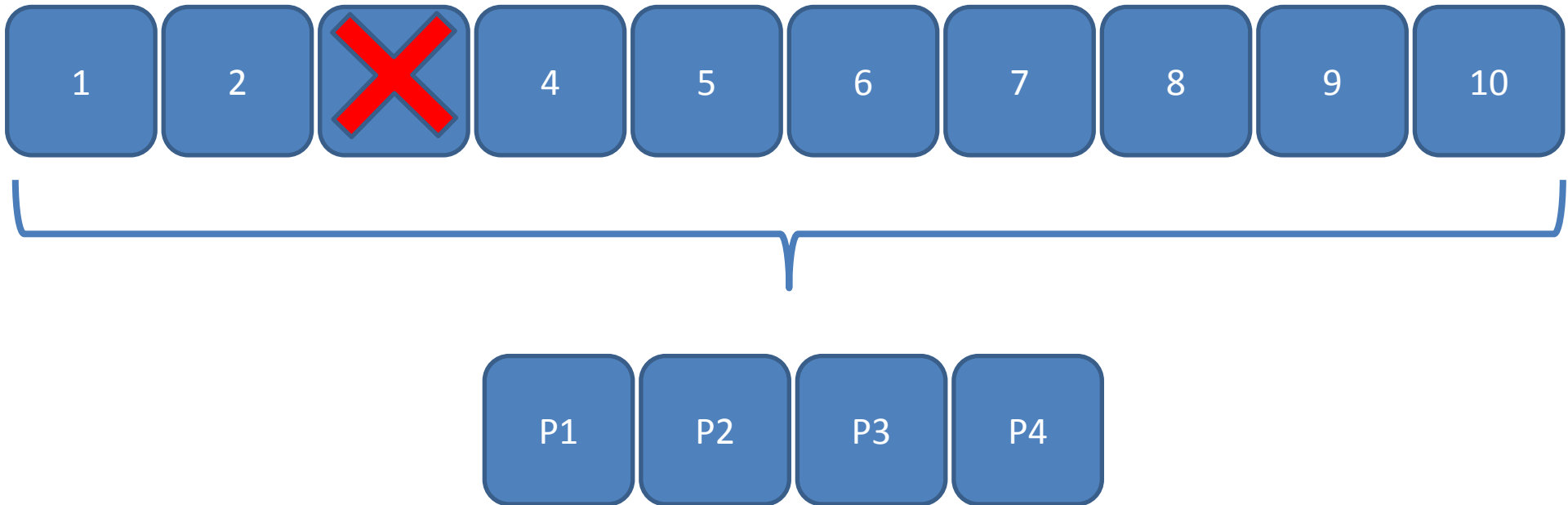
One failure results in one traffic unit

10:4 Code (Facebook)



Overhead: 40%

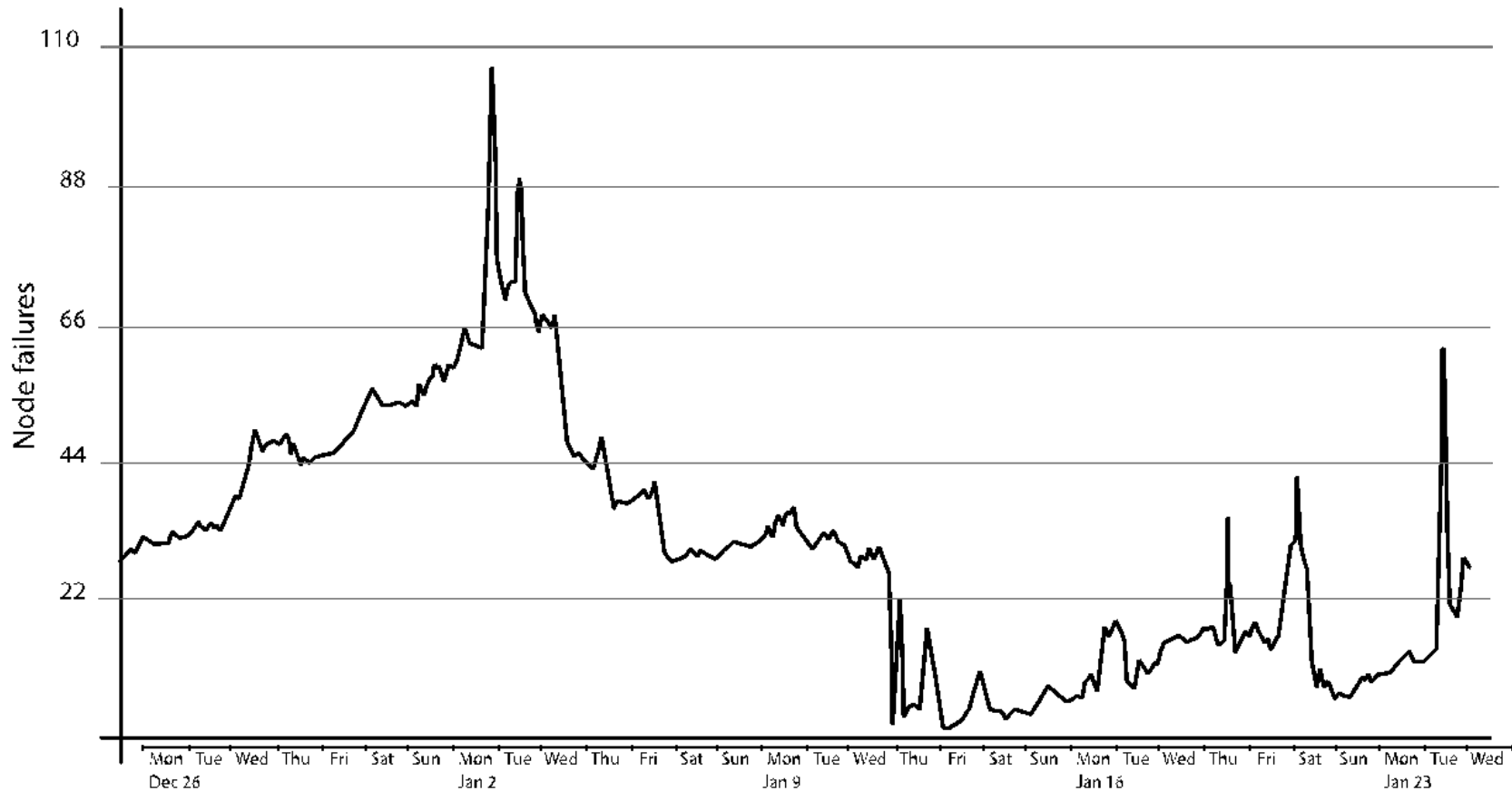
10:4 Code (Facebook)



Overhead: 40%

One failure results in 10 traffic units

Is repair frequent?



20 node failures * 15TB = 300TB

if 8% RS coded, 588TB network traffic/day. (average total network: 2PB/day)

~30% of network traffic is repair in a normal day