On Linear Field Size Access-Optimal MDS Convertible Codes

Myna Vajha

Indian Institute of Technology, Hyderabad

GMU Talk, June 30

Outline

- Erasure Codes and Distributed Storage.
 - MDS codes
- Code conversion problem
 - MDS convertible codes
 - Access and Bandwidth Cost
 - Access Optimal Constructions
 - Bandwidth Optimal Constructions
- Node repair problem (Overview)

How to save yourself from Data Loss ?

Redundant Array of independent Disks (RAID)

- RAID 1: uses replication
- RAID 3: adds a parity
- Distributed Storage (Replication or Erasure Codes)



DAID 2

Image Souces: Wikipedia, vifx.co.nz, technissile.com

Erasure Code

Used to recover from data loss.

- In a (n, k) erasure code, k units of data are encoded to get n units of data.
 - Storage Overhead of such a code $\nu = \frac{n}{k} = 1 + \frac{r}{k}$
 - erasure tolerance (number of erasures (any pattern) that can be corrected)
- A Maximum Distance Separable (MDS) Code provides reliability against erasure of any r = n k units.
 - ▶ (6,4) MDS code below has storage overhead 1.5 and it can recover from 2 erasures.



Replication Code vs Erasure Code

- Google's GFS uses 3-replication code.
- Facebook's Hadoop uses [14, 10] Reed Solomon(RS) Code.



(6,2) 3-Replication

[4,2] Reed Solomon

Code	Storage O/h	Bandwidth ¹	Reliability
[6,2] 3-rep	3x	0.5	2
[4,2] RS	2x	1.0	2

¹As a fraction of k units

MDS codes: Properties

• Let G = [I P] be $(k \times n)$ generator matrix of [n, k] code i.e.,

$$[p_1 \ p_2 \ \cdots \ p_{n-k}] = [m_1 \ \cdots \ m_k] \underbrace{\begin{bmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n-k} \\ p_{2,1} & p_{2,2} & \cdots & p_{2,n-k} \\ \vdots & \vdots & \ddots & \vdots \\ p_{k,1} & p_{k,2} & \cdots & p_{k,n-k} \end{bmatrix}}_{P}$$

 Say there are (n − k) erasures out of which m of them are message erasures given by i₁, ..., i_m and available parities given by j₁, ..., j_m.

$$[p_{j_1} \ p_{j_2} \ \cdots \ p_{j_m}] = [m_{i_1} \ \cdots \ m_{i_m}] \underbrace{\begin{bmatrix} p_{i_1,j_1} & p_{i_1,j_2} & \cdots & p_{i_1,j_m} \\ p_{i_2,j_1} & p_{i_2,j_2} & \cdots & p_{i_2,j_m} \\ \vdots & \vdots & \ddots & \vdots \\ p_{i_m,j_1} & p_{i_m,j_2} & \cdots & p_{i_m,j_m} \end{bmatrix}}_{\text{sub-matrix of } P}$$

▶ Code is MDS iff any $m \times m$ sub-matrix of P is invertible for all $m \le \min(k, n-k)$

MDS Code: Reed Solomon Construction

- [*n*, *k*] RS code.
 - A message of k symbols can be thought of as coefficients of a polynomial with degree < k.</p>
 - Evaluate the polynomial at n distinct points. (message symbols need not be present directly in these n evaluations)



k = 2, the polynomial is a line. Two evaluations enough to recover a line k = 4, the polynomial is of degree 3. Four evaluations enough to recover it

- If you have one evaluation of a line, then are infinitely many options for the lines that pass through the point.
- Need k evaluations to recover any polynomial of degree < k.

The Code Conversion Problem

Joint work with

- Nikhil Krishnan (Assistant Professor, IIT Palakkad)
- Vinayak Ramkumar (Postdoc, TU Munich)
- Xiangliang Kong, (Postdoc, Tel Aviv University)

The Code Conversion Problem: Motivation



As the disks age, it makes sense to

- Kadekodi et. al. "Cluster storage systems gotta have HeART: improving storage efficiency by exploiting disk-reliability heterogeneity". USENIX FAST, 2019.
- Kadekodi et. al, "PACEMAKER: Avoiding HeART attacks in storage clusters with disk-adaptive redundancy", USENIX OSDI 2020

Convertible Codes: Framework

- Convertible codes a framework introduced by Maturana and Rashmi [TIT 2022] to study the code conversion.
- Goal: To be able to effectively change from any [n' = k' + r', k'] initial code to a [n^F = k^F + r^F, k^F] final code.
 - Multiple initial codewords converted to multiple final codewords



► For $M = \text{lcm}(k^{I}, k^{F})$, $\lambda^{I} = \frac{M}{k^{I}}$ initial codewords get converted to $\lambda^{F} = \frac{M}{k^{F}}$ final codewords. Number of message symbols across initial and final codewords is M.

MDS Convertible Codes: Access cost

- MDS convertible codes are convertible codes where the initial and final codes are Maximum Distance Separable (MDS) codes.
- Access cost: number of symbols read from initial codewords to construct the final codewords plus the number of symbols written.
- We will focus on merge regime where $k^F = \lambda k'$.
 - default approach: download λk^{l} symbols to construct the parity symbols of the final code.
 - Tight lower bound on access cost (Maturana and Rashmi, TIT 22)

access cost
$$\geq \begin{cases} \lambda r^F + r^F & r^F \leq \min(r^I, k^I) \\ \lambda k^I + r^F & \text{otherwise} \end{cases}$$

- Assume $r^F \leq \min(r^I, k^I)$ the non-trivial case
- Split regime when $k' = \lambda^F k^F$

MDS Convertible Codes: An Example

Let G['] = [I_k, P[']], G^F = [I_k, P^F] be generator matrices of initial and final codes respectively. P['], P^F are of sizes k['] × r['] and λk['] × r^F respectively.

• Example:
$$k^{l} = 3, \lambda = 2, k^{F} = 6$$
 and $r^{l} = r^{F} = 2$

$${\cal P}' = \left[egin{array}{ccc} 1 & 1 \ heta_1 & heta_2 \ heta_1^2 & heta_2^2 \end{array}
ight], {\cal P}^{{\scriptscriptstyle F}} = \left[egin{array}{ccc} 1 & 1 \ heta_1 & heta_2 \ heta_1^2 & heta_2^2 \ heta_1^2 & heta_2^2 \ heta_1^3 & heta_2^3 \ heta_1^4 & heta_2^4 \ heta_1^5 & heta_2^5 \end{array}
ight]$$

• Let p_1^1, p_2^1 and p_1^2, p_2^2 be the parities from two initial codewords.

$$\begin{bmatrix} p_1^{I} & p_2^{I} \end{bmatrix} = \begin{bmatrix} m_1^{I} & m_2^{I} & m_3^{I} \end{bmatrix} P_I \implies p_i^{I} = m_1^{I} + m_2^{I} \theta_i + m_3^{I} \theta_i^{2} \\ \begin{bmatrix} p_1^{F} & p_2^{F} \end{bmatrix} = \begin{bmatrix} m_1^{1} & m_1^{1} & m_1^{1} & m_1^{2} & m_2^{2} & m_3^{2} \end{bmatrix} P_F$$
$$p_1^{F} = p_1^{1} + \theta_1^{3} p_1^{2} \\ p_2^{F} = p_2^{1} + \theta_2^{3} p_2^{2}$$

• Can do conversion in this case by accessing $\lambda r^F = 4 < \lambda k^I = 6$ symbols

Access Optimality from Block-Reconstructable Property

$$P^{F} = \begin{bmatrix} P^{F,1} \\ P^{F,2} \\ \vdots \\ P_{F,\lambda} \end{bmatrix} \text{ where } P^{F,\ell} \text{ is } k^{I} \times r^{F} \text{ matrix.}$$

- *P^F* is said to be *r^F*-block reconstructable from *P^I* if for each *l* ∈ [λ] there exists *r^F* columns of *P^I* that span the columns of *P^F*.
- MDS convertible code is access-optimal if P^F is r^F -block reconstructable from P^I
 - r^F final parities can be constructed by accessing exactly λr^I parity symbols.

Per-Symbol Access Optimality

- Recovery of each final parity symbol uses exactly λ initial parities with each parity belonging to an initial codeword.
- P^F is said to be parallel-block reconstructable from P^I if for each $\ell \in [\lambda]$ there exist r^F columns of P^I that are exactly equal to or scaling of columns seen in $P^{F,\ell}$.

parallel-block reconstructable \rightarrow block-reconstructable

• Helps non-central conversion setting where the new node downloads the required data to reconstruct the corresponding final parity.

Earlier Approaches

Vandermonde construction

$$P' = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \theta_1 & \theta_2 & \cdots & \theta_{r'} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_1^{k'-1} & \theta_2^{k'-1} & \cdots & \theta_{r'}^{k'-1} \end{bmatrix}, P^F = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \theta_1 & \theta_2 & \cdots & \theta_{r'} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_1^{\lambda k'-1} & \theta_2^{\lambda k'-1} & \cdots & \theta_{r'}^{\lambda k'-1} \end{bmatrix}$$

▶ Need q to be large $(O(2^{(n^F)^3}))$ to guarantee that P^F and P^I are super-regular

2 Hankel matrix based convertible codes

- Super-regular property guaranteed for sub-matrices of Hankel matrices.
- ► Constructions limited to parameters r^F ≤ r^I − λ + 1 with linear field size for fixed number of parities

Maturana and Rashmi, "Convertible Codes: New Class of Codes for Efficient Conversion of Coded Data in Distributed Storage", ITCS 2020

Earlier Approaches: Polynomial Construction

$$H_{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_{1} & x_{2} & \cdots & x_{m} \\ x_{1}^{2} & x_{2}^{2} & \cdots & x_{m}^{2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1}^{\prime - 1} & x_{2}^{\prime - 1} & \cdots & x_{m}^{\prime - 1} \end{bmatrix}, \text{ where } X = \{x_{1}, \cdots, x_{m}\}$$

• Initial code: An (n', k') MDS code defined by parity check matrix: $H' = [H_{X_1} H_Y] \implies -H_{X_1} \underline{m}^j = H_Y \underline{p}^j$

• Final code: An (n', k') MDS code defined by parity check matrix:

$$H^{F} = [H_{X_{1}} H_{X_{2}} \cdots H_{X_{\lambda}} H_{Y}] \implies -\sum_{j=1}^{\lambda} H_{X_{j}} \underline{m}^{j} = H_{Y} \underline{p}^{F}$$

• Suppose $H_{X_j} = D_j H_{X_1}$ where D_j is $r' \times r'$ diagonal matrix

$$\underline{p}^{F} = -H_{Y}^{-1} \sum_{j=1}^{\lambda} H_{X_{j}} \underline{m}^{j} = -H_{Y}^{-1} \sum_{j=1}^{\lambda} D_{j} H_{X_{1}} \underline{m}^{j} = -H_{Y}^{-1} \sum_{j=1}^{\lambda} D_{j} \underline{p}^{j}$$

⇒ access-optimal MDS code. Pick $X_j = \gamma^j X_1$ such that X_i 's, Y are disjoint. Can add a column to H_Y such that $q \ge k^F + r^I - 1$ is sufficient

Xiangliang Kong, "Locally Repairable Convertible Codes With Optimal Access Costs", IEEE Trans. in Info. Theory 2024.
 Krishnan et. al. "On Low Field Size MDS Convertible Codes", IEEE ISIT 2025

Our Approach

- P^{I} and P^{F} are Cauchy matrices.
- Cauchy matrix C(X, Y) where $X = \{x_1, \dots, x_k\}$, $Y = \{y_1, \dots, y_r\}$:

$$C(X,Y) = \begin{bmatrix} \frac{1}{x_1 - y_1} & \frac{1}{x_1 - y_2} & \cdots & \frac{1}{x_1 - y_r} \\ \frac{1}{x_2 - y_1} & \frac{1}{x_2 - y_2} & \cdots & \frac{1}{x_2 - y_r} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{x_k - y_1} & \frac{1}{x_k - y_2} & \cdots & \frac{1}{x_k - y_r} \end{bmatrix}$$

- Cauchy matrices are super-regular.
- Goal: Design X, Y and X_1 of cardinalities k^F , r^I , k^I such that:

 $P^F = C(X, Y)$ is parallel-block reconstructable from $P^I = C(X_1, Y)$

- Enough to look at $r^F = r^I = r$.
- For $r^F \leq r^I$, we do not convert $(r^I r^F)$ nodes

Sub-Group Based Construction: Using Multiplicative Sub-Group of \mathbb{F}_q^*

- Let \mathbb{F}_q be such that $r \mid (q-1)$ and $\nu = \frac{q-1}{r}$, $k' \leq \nu 1$
- Let α be the primitive element of F_q and γ = α^ν.

 $Y = \{\gamma, \gamma^2, \cdots, \gamma^r = 1\}, \quad X_1 = \{\alpha, \alpha^2, \cdots, \alpha^{k'}\}, \quad X_\ell = \gamma^{\ell-1}X_1, X = \cup_{\ell=1}^{\lambda}X_\ell$

• $X_1, \cdots, X_{\lambda}, Y$ are disjoint if $\lambda \leq r$

• $P^F = C(X, Y)$ is parallel-block reconstructable from $P^I = C(X_1, Y)$

$$P^{F,\ell}(i,j) = \frac{1}{\gamma^{\ell-1}\alpha^i - \gamma^j} = \frac{\gamma^{-(\ell-1)}}{\alpha^i - \gamma^{j-\ell+1}}$$
$$= \gamma^{-(\ell-1)}P'(i,j')$$

where $j' \in [r]$ such that $\gamma^{j'} = \gamma^{j-\ell+1}$.

• *j*-th column of $P^{F,\ell}$ is scaling of *j'*-th column of P^I

Sub-Group Based Construction: Using Multiplicative Sub-Group of \mathbb{F}_q^*

- Modification 1: $\lambda \leq (r-1)$
 - ▶ Let \mathbb{F}_q be such that $(r-1) \mid (q-1)$ and $\nu = \frac{q-1}{r-1}$, $k^l \leq \nu 1$, $\gamma = \alpha^{\nu}$.

$$Y = \{\mathbf{0}, \gamma, \cdots, \gamma^{r-1} = 1\}$$

►
$$P^F = C(X, Y)$$
 is parallel-block reconstructable from $P^I = C(X_1, Y)$
★ $P^{F,\ell}(i, 1) = \frac{1}{\gamma^{\ell-1}\alpha^i} = \gamma^{-(\ell-1)}P^I(i, 1)$

- Modification 2: $\lambda \leq (r-2)$
 - Append an all-one column to Cauchy matrix is still super-regular matrix (Roth and G. Seroussi, TIT 85)
 - \mathbb{F}_q be such that $(r-2) \mid (q-1)$ and $\nu = \frac{q-1}{r-2}$, $k' \leq \nu 1$, $\gamma = \alpha^{\nu}$.

$$Y = \{\mathbf{0}, \gamma, \cdots, \gamma^{r-2} = 1\},\$$

▶
$$P^F = \begin{bmatrix} 1 & C(X, Y) \end{bmatrix}$$
 is parallel-block reconstructable from $P^I = \begin{bmatrix} 1 & C(X_1, Y) \end{bmatrix}$

Sub-Group Based Construction: Using Multiplicative Sub-Group of \mathbb{F}_q^*

•
$$r' = r^F = 4, k' = 5, \lambda = 2 \implies k^F = 10, n^F = 14$$

• Let q = 13 (meets the MDS conjecture $q = n^{F} - 1$)

$$\begin{split} Y &= \{2^6 = 11, 2^{12} = 1, 0\}, X_1 = \{2, 2^2, 2^3, 2^4, 2^5\}, \\ &X_2 = \{2^6, 2^7, 2^8, 2^9, 2^{10}, 2^{11}\} \end{split}$$

• If
$$\lambda = r - 2$$
 and $(r - 2) \mid (q - 1)$ for
 $q = (k^{l} + 1)(r - 2) + 1 = n^{F} - 1$, then
MDS conjecture is met.

						_	
j	Γ	1	9	7	1	l	
į		9	8	10	1	l	
į		2	3	5	1		P^{I}
		7	10	9	1		
ļ		8	2	11	1	ľ	
		4	12	6	1	Γ	
		5	4	3	1		
	1	.0	11	8	1		
		3	6	4	1		
	1	.1	5	2	1		

Sub-Group Based Construction: Additive Sub-Group of \mathbb{F}_q

• Let
$$q = p^m$$
 then \mathbb{F}_q is isomorphic to $\{f(x) \in \mathbb{F}_p[x] \mid \deg(f(x)) < m\}$
• $r = p^u$ i.e., $r \mid q$ and $\nu = \frac{q}{r} = p^{m-u}$
 $Y = \{f(x) \in \mathbb{F}_p[x] \mid \deg(f(x)) < u\} = \{y_1(x), \dots, y_r(x)\}$
 $X_1 = \{f_1(x), \dots, f_k(x)\}$
 $\subseteq \{x^u f(x) \mid f(x) \in \mathbb{F}_p[x], \deg(f(x)) < m - u - 1\}$
 $X_\ell = y_\ell(x) + X_1$

• $P^F = C(X, Y)$ is parallel-block reconstructable from $P^I = C(X_1, Y)$

$$P^{F,\ell}(i,j) = \frac{1}{y_{\ell}(x) + f_i(x) - y_j(x)} = \frac{1}{f_i(x) - y_{j'}(x)}$$
$$= P^{I}(i,j')$$

where $j' \in [r]$ such that $y_{j'}(x) = y_j(x) - y_\ell(x)$.

• *j*-th column of $P^{F,\ell}$ is same as *j*'-th column of P^I

• Can similarly modify to add all one column.

Access Cost: Split regime

$$ext{access cost} \geq egin{cases} (\lambda-1)k^F + \min(r^F,k^F) + r^F & r^I \geq r^F \ \lambda k^F + r^F & ext{otherwise} \end{cases}$$

•
$$(n' = k' + r^F, k' = \lambda k^F, n^F = k^F + r^F, k^F)$$

• Consider any initial MDS code described by generator matrix G' = [I P']

- Let $G^F = [I P^F]$ be generator matrix of final code where $P^F = P^I(1:k^F, 1:r^F)$
- Let m^1, \cdots, m^{λ} be message symbols of the λ final codewords
- Download m^2, \dots, m^{λ} and compute $m^2 p^F, \dots, m^{\lambda} p^F$
- Download the r^F initial parities given by:

$$[m_1 \cdots m_{\lambda}]P^{I} = m_1P^{F} + \underbrace{[m_2 \cdots m_{\lambda}]}_{known}P^{I}(k^{F}+1:k^{I},1:r^{F})$$

recover $m_1 P^F$ from above.

Central vs Cooperative Conversion



• Write cost is $r^F \alpha$ in the central scheme

• Can be made to be much smaller in the cooperative scheme.

Maturana and Rashmi: "Bandwidth Cost of Code Conversions in Distributed Storage: Fundamental Limit and Optimal Constructions", Trans. in Info. Theory 2023

Bandwidth Cost of Code Conversion: Merge Regime

- α is the number of symbols stored at each node.
 - Example: $k^{I} = 4, r^{I} = 1, r^{F} = 3$. Two initial codewords shown below. Each code symbol is a vector with $\alpha = r^{F} = 3$ symbols in F_{q} . (Maturana and Rashmi, TIT 2023)

$m_{0,0}^1$	$m_{0,1}^1$	m ¹ _{0,2}	$m_{0,0}^2$	$m_{0,1}^2$	$m_{0,2}^2$
$m_{1,0}^1$	$m_{1,1}^1$	$m_{1,2}^1$	$m_{1,0}^2$	$m_{1,1}^2$	$m_{1,2}^2$
$m_{2,0}^1$	$m_{2,1}^1$	$m_{2,2}^1$	$m_{2,0}^2$	$m_{2,1}^2$	$m_{2,2}^2$
$m_{3,0}^1$	$m_{3,1}^1$	m ¹ _{3,2}	$m_{3,0}^2$	$m_{3,1}^2$	$m_{3,2}^2$
$m_0^1 p^{(1,0)}$	$m_1^1 p^{(1,0)} + m_0^1 p^{(1,1)}$	$m_2^1 p^{(1,0)} + m_0^1 p^{(1,2)}$	$m_0^2 p^{(1,0)}$	$m_1^2 p^{(1,0)} + m_0^2 p^{(1,1)}$	$m_2^2 p^{(1,0)} + m_0^2 p^{(1,2)}$

$m_{0,0}^1$	$m_{0,1}^1$	$m_{0,2}^1$
$m_{1,0}^1$	$m_{1,1}^1$	$m_{1,2}^1$
$m_{2,0}^1$	$m_{2,1}^1$	$m_{2,2}^1$
$m_{3,0}^1$	$m_{3,1}^1$	$m_{3,2}^1$
$m_{0,0}^2$	$m_{0,1}^2$	$m_{0,2}^2$
$m_{1,0}^2$	$m_{1,1}^2$	$m_{1,2}^2$
$m_{2,0}^2$	$m_{2,1}^2$	$m_{2,2}^2$
$m_{3,0}^2$	$m_{3,1}^2$	$m_{3,2}^2$
$m_0^1 p^{(1,0)} + m_0^2 p^{(2,0)}$	$m_1^1 p^{(1,0)} + m_1^2 p^{(2,0)}$	$m_2^1 p^{(1,0)} + m_2^2 p^{(2,0)}$
$m_0^1 p^{(1,1)} + m_0^2 p^{(2,1)}$	$m_1^1 p^{(1,1)} + m_1^2 p^{(2,1)}$	$m_2^1 p^{(1,1)} + m_2^2 p^{(2,1)}$
$m_0^1 p^{(1,2)} + m_0^2 p^{(2,2)}$	$m_1^1 p^{(1,2)} + m_1^2 p^{(2,2)}$	$m_2^1 p^{(1,2)} + m_2^2 p^{(2,2)}$

Maturana and Rashmi: "Bandwidth Cost of Code Conversions in Distributed Storage: Fundamental Limit and Optimal Constructions", Trans. in Info. Theory 2023

Conclusions and Open Questions

- Per-symbol access-optimal constructions for $\lambda \leq r$ with low field size.
 - Can this parameter restriction be removed without increasing the field size ?
- What are the fundamental bandwidth limits for code conversion in split regime ?
- What are the fundamental bandwidth limits of co-operative conversion ?
 - We have some preliminary schemes for the cooperative case that have smaller bandwidth requirements compared to central conversion bandwidth

The Node Repair Problem

Erasure Codes and Node Failures



- A median of 50 nodes are unavailable per day.
- 98% of the failures are single node failures.
- A median of 180TB of network traffic per day is generated in order to reconstruct the RS coded data corresponding to unavailable machines.
- Thus there is a strong need for erasure codes that can efficiently recover from single-node failures.

Image courtesy: Rashmi et al.: "A Solution to the Network Challenges of Data Recovery in Erasure-coded Distributed Storage Systems: A Study on the Facebook Warehouse Cluster," USENIX Hotstorage, 2013.

Conventional Node Repair of an RS Code



In the example (14, 10) RS code,

the amount of data downloaded to repair 100MB of data equals 1GB.

clearly, there is room for improvement...

Things that matter.

- Low Storage Overhead. ✓
- High Tolerance for erasures. \checkmark
- Fast Repair (Single Node Failure). ??
 - Computation efficiency.
 - Speed in procuring repair data.
 - ★ low repair traffic.
 - ★ smaller disk read latencies.





Image source: sine.co, cartoonstock.com

Regenerating Codes

Parameters: $((n, k, d), (\alpha, \beta), B, \mathbb{F}_q)$



- Data (of size *B*) can be recovered by connecting to any *k* of *n* nodes
- A failed node can be repaired by connecting to any *d* nodes, downloading β symbols from each node; (*d*β << file size *B*)

A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems", IEEE Transactions on Information Theory, 2010

Minimum Storage Regenerating Code





- Size of failed node's contents: 100MB
- 2 RS repair BW: 1 GB
- MSR Repair BW: 325 MB

4-way Optimality of Clay code





Image courtesy: denverpost.com

Vajha et. al., Clay Codes: Moulding MDS Codes to Yield an MSR Code, USENIX FAST 2018.

Ceph: Contributions



A popular opensource distributed storage system used by CERN, Flipkart, Cisco etc



"Us (+Vinayak) pitching Clay codes to Ceph in April 2017"

We have introduced Clay code as erasure code plugin. It is part of Ceph's Nautilus release (March 2019). As part of this we also introduced support for vector codes in Ceph.

Vajha et. al., Clay Codes: Moulding MDS Codes to Yield an MSR Code, USENIX FAST 2018.

Clay Code Summary

- The open-source implementation of Clay code in Ceph is for any (n, k, d) parameters.
- In comparison to (20, 16) RS code, for Workloads with large sized objects (64MB), the Clay code (20, 16, 19):
 - resulted in repair time reduction by 3X.
 - Improved degraded read and write performance by 27.17% and 106.68% respectively.

Vajha et. al., Clay Codes: Moulding MDS Codes to Yield an MSR Code, USENIX FAST 2018.

Open Questions

- Sub-packetization level prohibitive in realizing the benefits of MSR codes in practice.
- Given a fixed sub-packetization level, what is the minimal repair bandwidth that is required ?

Appendix

Moulding an MDS Code to Yield the Clay Code

(n = 4, k = 2) MDS code with optimal repair of systematic nodes, $\alpha = 2$



Code symbols of [4,2] MDS code.

Layer two such units

Uncoupled code still needs 4 symbols during recovery of single node (containing 2 symbols).

Moulding an MDS Code to Yield the Clay Code

(n = 4, k = 2) MDS code with optimal repair of systematic nodes, $\alpha = 2$



Code symbols of [4,2] MDS code.

Layer two such units

Uncoupled code still needs 4 symbols during recovery of single node (containing 2 symbols).

Moulding an MDS Code to Yield the Clay Code

(n = 4, k = 2) MDS code with optimal repair of systematic nodes, $\alpha = 2$



- Uncoupled code has 2 planes, where each plane corresponds to an [4, 2] MDS code
- Coupled code symbols are obtained by:
 - Copying symbols with red dots
 - Pair of yellow symbols {C, C*} are obtained by transformation

• Note that recovery of any failed node in Uncoupled code requires 4 symbols







symbols that are computable in uncoupled cube









Clay Code

(n = 4, k = 2, d = 3) MSR code with all node optimal repair



Coupled Code



• The same construction extends to any (n, k, d)