

On Codes for Optimal Rebuilding Access

Zhiying Wang*, Itzhak Tamo*[†], and Jehoshua Bruck*

*Electrical Engineering Department, California Institute of Technology, Pasadena, CA 91125, USA

[†]Electrical and Computer Engineering, Ben-Gurion University of the Negev, Beer Sheva 84105, Israel
 {zhiying, tamo, bruck}@caltech.edu

Abstract—MDS (maximum distance separable) array codes are widely used in storage systems due to their computationally efficient encoding and decoding procedures. An MDS code with r redundancy nodes can correct any r erasures by accessing (reading) all the remaining information in both the systematic nodes and the parity (redundancy) nodes. However, in practice, a single erasure is the most likely failure event; hence, a natural question is how much information do we need to access in order to rebuild a single storage node? We define the *rebuilding ratio* as the fraction of remaining information accessed during the rebuilding of a single erasure. In our previous work we showed that the optimal rebuilding ratio of $1/r$ is achievable (using our newly constructed array codes) for the rebuilding of any systematic node, however, all the information needs to be accessed for the rebuilding of the parity nodes. Namely, constructing array codes with a rebuilding ratio of $1/r$ was left as an open problem. In this paper, we solve this open problem and present array codes that achieve the lower bound of $1/r$ for rebuilding any single systematic or parity node.

I. INTRODUCTION

MDS (maximum distance separable) array codes are a family of erasure-correcting codes used extensively as the basis for RAID storage systems. An array code consists of a 2-D array where each column can be considered as a disk. We will use the term column, node, or disk interchangeably. A code with r parity (redundancy) nodes is MDS if and only if it can recover from any r erasures. EVENODD [2] and RDP [5] are examples of MDS array codes with two redundancies. In this paper, we only consider systematic codes, namely, the information is stored exclusively in the first k nodes, and the parities are stored exclusively in the last r nodes.

In order to correct r erasures, it is obvious that one has to access (or read) the information in all the surviving nodes. However, in practice it is more likely to encounter a single erasure rather than r erasures. So a natural question is: How much information do we need to access when rebuilding a single erasure? Do we have to access all the surviving information? We define the *rebuilding ratio* as the ratio of accessed information to the remaining information in case of a single erasure. For example, it is easy to check that for the code in Figure 1, if any two columns are erased, we can still recover all the information, namely, it is an MDS code. However, if column C_1 is erased, it can be rebuilt by accessing $a_{0,2}, a_{1,2}$ from column C_2 , r_0, r_1 from column C_3 , and z_0, z_1 from column C_4 , as follows:

$$\begin{aligned} a_{0,1} &= 2a_{0,2} + r_0 \\ a_{1,1} &= 2a_{1,2} + r_1 \end{aligned}$$

Systematic nodes		Parity nodes	
C_1	C_2	C_3	C_4
$a_{0,1}$	$a_{0,2}$	$r_0 = a_{0,1} + a_{0,2}$	$z_0 = a_{2,1} + a_{1,2}$
$a_{1,1}$	$a_{1,2}$	$r_1 = a_{1,1} + a_{1,2}$	$z_1 = a_{3,1} + 2a_{0,2}$
$a_{2,1}$	$a_{2,2}$	$r_2 = a_{2,1} + a_{2,2}$	$z_2 = 2a_{0,1} + 2a_{3,2}$
$a_{3,1}$	$a_{3,2}$	$r_3 = a_{3,1} + a_{3,2}$	$z_3 = 2a_{1,1} + a_{2,2}$

Figure 1. An MDS array code with two systematic and two parity nodes. All the elements are in finite field F_3 . The first parity column C_3 is the row sum and the second parity column C_4 is generated by the zigzags. For example, zigzag z_0 contains the elements $a_{i,j}$ that satisfy $f_j^1(i) = 0$.

$$\begin{aligned} a_{2,1} &= 2a_{1,2} + z_0 \\ a_{3,1} &= a_{0,2} + z_1 \end{aligned}$$

Here all elements are in finite field F_3 . Hence, by accessing only half of the remaining information, the erased node can be rebuilt. Details on this new code will be discussed in Section II.

A related problem called *repair bandwidth* was first proposed in [6]. The paradigm there is that one can access the entire information and perform computations within each node, and the question is how much information is *transmitted* for rebuilding? A lower bound on the repair bandwidth was given in [6]. When a single erasure occurs and all the remaining nodes are accessible, the lower bound for the bandwidth is $\frac{1}{r}$. Recently, a number of codes were designed to achieve the bandwidth lower bound. When the number of parity nodes is larger than that of the systematic nodes, explicit code constructions were given in [8]–[10]. For all cases, [4], [11] achieved the lower bound asymptotically.

It is clear that a lower bound on the repair bandwidth is also a lower bound on the rebuilding ratio. In [12] we presented an explicit construction of MDS array codes that achieve the lower bound on the ratio for rebuilding any *systematic node*. A similar code construction was given in [3]. Also in [7] a similar code with 2 parities was proposed - it has optimal repair bandwidth for any single erasure.

The main contribution of this paper is an explicit construction of MDS array codes with r parity nodes, that achieves

the lower bound $1/r$ for rebuilding *any systematic or parity node*. The rebuilding of a single erasure has an efficient implementation as computations within nodes are not required. Moreover, our codes have simple encoding and decoding procedures - when $r = 2$ and $r = 3$, the codes require finite-field sizes of 3 and 4, respectively.

The rest of the paper is organized as follows. Section II introduces the rebuilding ratio problem for MDS array codes and reviews the code construction in [12]. Section III describes the construction of our codes with optimal rebuilding ratio. Finally, the paper is summarized in Section IV.

II. REBUILDING RATIO PROBLEM

In this section we formally define the rebuilding ratio problem and review the code construction in [12]. We then prove that the construction can be made an MDS code, in fact, this will be the basis for proving that our newly proposed construction which is described in Section III is also an MDS code.

We first define the framework of a systematic MDS array code. Let $A = (a_{i,j})$ be an information array of size $p \times q$. A column is also called a node, and an entry is called an element. Each of the q columns is a systematic node in the code. We add r parity columns to this array on the right, such that from any q columns, we can recover the entire information. In [12], it was shown that if each information element is protected by exactly r parity elements, then each parity node corresponds to q permutations acting on $[0, p-1]$. More specifically, suppose the permutations are f_1, f_2, \dots, f_q . Then the t -th element in this parity node is a linear combination of all elements $a_{i,j}$ such that $f_j(i) = t$. The set of information elements contained in this linear combination is called a *zigzag set*. For the t -th element in the l -th parity, $t \in [0, p-1], l \in [0, r-1]$, denote by f_1^l, \dots, f_q^l the set of associated permutations, and Z_t^l the zigzag set.

The ordering of the elements in each node can be arbitrary, hence, we can assume that the first parity node is always a linear combination of each row (corresponding to identity permutations). Figure 1 is an example of such codes. The first parity C_3 corresponds to identity permutations. The second parity C_4 corresponds to the permutations

$$\begin{aligned} f_1^1 &= (2, 3, 0, 1), \\ f_2^1 &= (1, 0, 3, 2). \end{aligned}$$

For a given MDS code with parameters q, r , we ask what is the accessed fraction in order to rebuild a single node (in the average case)? Hence, the *rebuilding ratio* of a code is:

$$R = \frac{\sum_{i=1}^{q+r} (\# \text{ accessed elements to rebuild node } i)}{(q+r)(\# \text{ remaining elements})}.$$

When a systematic node is erased, we rebuild each unknown element by one of the parity nodes. That is, we access one parity element containing the unknown, and access all the elements in the corresponding zigzag set except the unknown. In order to lower the number accesses, we would like to

find (i) good permutations such that the accessed zigzag sets intersect as much as possible, and (ii) proper coefficients in the linear combinations such that the code is MDS. For example, in Figure 1, in order to rebuild column C_1 , we access the zigzag sets $A = \{Z_0^0, Z_1^0\}, B = \{Z_0^1, Z_1^1\}$, corresponding to parities $\{r_0, r_1\}, \{z_0, z_1\}$. The surviving elements in A and in B are identical, i.e., $\{a_{0,2}, a_{1,2}\}$, therefore, only $1/2$ of the elements are accessed. Besides, the coefficients $\{1, 2\}$ in the parity linear combinations guarantee that any two nodes are sufficient to recover all the information. Hence the code is MDS.

Next we review the construction with optimal rebuilding for systematic nodes that was presented in [12]. The idea in the code construction was to form permutations based on r -ary vectors.

Let e_1, e_2, \dots, e_k be the standard vector basis of \mathbb{Z}_r^k . We will use x to represent both an integer in $[0, r^k - 1]$ and its r -ary expansion (the r -ary vector of length k). It will be clear from the context which meaning is used. All the calculations are done over \mathbb{Z}_r .

Construction 1 Let the information array be of size $r^k \times k$. Define permutation f_j^l on $[0, r^k - 1]$ as $f_j^l(x) = x + le_j$, $j \in [1, k], l \in [0, r-1]$. For $t \in [0, r^k - 1]$, we define the zigzag set Z_t^l in parity node l as the elements $a_{i,j}$ such that their coordinates satisfy $f_j^l(i) = t$. Let $Y_j = \{x \in [0, r^k - 1] : x \cdot e_j = 0\}$. Rebuild column j by accessing rows Y_j in all remaining columns.

Theorem 1 Construction 1 has optimal ratio $1/r$ for rebuilding any systematic node [12].

Figure 1 is an example of Construction 1. As mentioned before, only $1/2$ of the information is accessed in order to rebuild C_1 . The accessed elements are in rows $Y_1 = \{x \in [0, 3] : x \cdot e_1 = 0\} = \{0, 1\}$.

Next, we show that by assigning the coefficients in the parities properly, the code is MDS. Let $P_j = (a_{i,l})$ be the permutation matrix corresponding to $f_j = f_j^1$, namely, $a_{i,l} = 1$ if $l + e_j = i$, and $a_{i,l} = 0$ otherwise. Assigning the coefficients is the same as modifying $a_{i,l} = 1$ to other non-zero values. When $r = 2, 3$, modify $a_{i,l} = 1$ to $a_{i,l} = c$, if $l \cdot \sum_{t=1}^j e_t = 0$, where c is a primitive element of F_3, F_4 , respectively. The above assignment will make the code MDS for $r = 2, 3$ [12]. For example, the coefficients in Figure 1 is assigned in this way.

When $r \geq 4$, modify all $a_{i,l} = 1$ to $a_{i,l} = \lambda_j$, for some λ_j in a finite field F . Let the generator matrix of the code be

$$G' = \begin{bmatrix} I & & & \\ & \ddots & & \\ & & I & \\ I & \dots & I & \\ P_1^1 & \dots & P_k^1 & \\ \vdots & & \vdots & \\ P_1^{r-1} & \dots & P_k^{r-1} & \end{bmatrix}.$$

The following theorem shows that under this assignment the code can be MDS.

Theorem 2 (1) Construction 1 can be made an MDS code for a large enough finite field.

(2) When $r = 2, 3$, field of size 3 and 4 is sufficient to make the code MDS.

Proof: Part (2) was given in [12]. We only prove part (1). An MDS code means that it can recover any r erasures. Suppose t systematic nodes and $r - t$ parity nodes are erased, $1 \leq t \leq r$. Thus suppose we delete from G' the systematic rows $\{j_1, j_2, \dots, j_t\}$ and the remaining parity nodes are $\{i_1, i_2, \dots, i_t\}$. Then the following $t \times t$ block matrix should be invertible:

$$G = \begin{bmatrix} P_{j_1}^{i_1} & \cdots & P_{j_t}^{i_1} \\ \vdots & & \vdots \\ P_{j_1}^{i_t} & \cdots & P_{j_t}^{i_t} \end{bmatrix} \quad (1)$$

Its determinant $\det(G)$ is a polynomial with indeterminates $\lambda_{j_1}, \dots, \lambda_{j_t}$. All terms have highest degree $r^k(i_1 + \dots + i_t)$. One term with highest degree is $\prod_{s=1}^t \lambda_{j_s}^{i_s r^k}$ with non-zero coefficient 1 or -1 . So $\det(G)$ is a non-zero polynomial. Up to now we only showed one possible case of erasures. For any r erasures, we can find the corresponding non-zero polynomial. The product of all these polynomials is again a non-zero polynomial. Hence by [1] for a large enough field there exist assignments of $\{\lambda_j\}$ such that the polynomial is not 0. Then each G is invertible, and the code is MDS. ■

III. CODE CONSTRUCTION

The code in [12] has optimal rebuilding for systematic nodes. However, in order to rebuild a parity node, one has to access all the information elements. In this section we construct MDS array codes with optimal rebuilding ratio for rebuilding both the systematic and the parity nodes. The code has $k - 1$ systematic nodes and r parities nodes, for any k, r .

Consider the permutation $f_j = f_j^1$ in Construction 1. It is clear that f_j is a permutation of order r , i.e., f_j^r is the identity permutation. For $i \in [0, r - 1]$, define X_i as the set of vectors of weight i , namely, $X_i = \{v \in \mathbb{Z}_r^k : v \cdot (1, \dots, 1) = i\}$. X_0 is a subgroup of \mathbb{Z}_r^k and $X_i = X_0 + ie_k$ is its coset, where $e_k = (0, \dots, 0, 1)$. Assume the elements in X_i are ordered, $i \in [0, r - 1]$, and the ordering is

$$\begin{aligned} X_0 &= (v_1, \dots, v_{r^{k-1}}), \\ X_i &= (v_1 + ie_k, \dots, v_{r^{k-1}} + ie_k). \end{aligned}$$

Since the ordering of the elements in each column does not matter, we can reorder them as $(X_0, X_1, \dots, X_{r-1})$, with each X_i ordered as above. One can check that $f_j(X_i) = X_{i+1}$, where the subscript is added mod r . So the matrix P_j can be

$$\begin{aligned} A^0 &= \begin{bmatrix} \underline{I} & & \\ \underline{p} & \alpha p^2 & \\ \underline{p^2} & & p \end{bmatrix} \\ A^1 &= \begin{bmatrix} \underline{p} & p^2 & \\ & I & \\ & p & \alpha p^2 \end{bmatrix} \\ A^2 &= \begin{bmatrix} \alpha p^2 & & p \\ & p & p^2 \\ & & I \end{bmatrix} \end{aligned}$$

Figure 2. Parity matrices A^i for $r = 2$ (left) and $r = 3$ (right) parities. When the first parity node is erased, the underlined elements are accessed from systematic nodes. The remaining unknown elements are recovered by the shaded elements from parity nodes.

written as

$$P_j = \begin{matrix} X_0 \\ X_1 \\ \vdots \\ X_{r-1} \end{matrix} \begin{pmatrix} X_0 & X_1 & \cdots & X_{r-1} \\ & p_j & & p_j \\ & & \ddots & \\ & & & p_i \end{pmatrix}, \quad (2)$$

where p_j corresponds to the mapping of $f_j : X_i \mapsto X_{i+1}$. In particular, if p_j is viewed as a permutation acting on X_0 , then for $x \in X_0$,

$$p_j(x) = x + e_j - e_k.$$

When $r = 2, 3$, modify the 1 entries of p_i into c if its corresponding column l satisfies $l \cdot \sum_{t=1}^j e_t = 0$. Here c is an primitive element in F_3, F_4 . When $r \geq 4$, modify 1 entries into λ_j .

In the following, we will use blocks the same as single elements. When referring to row or column indices, we mean block row or column indices. We refer to p_j as a small block, and the corresponding block row or column as a small block row or column. And P_j is called a big block with big block row or column. Moreover, we assume the elements in each column are in order (X_0, \dots, X_{r-1}) .

Construction 2 Suppose the information array is of size $r^k \times (k - 1)$. For $j \in [1, k - 1]$, define a big block matrix

$$A_j^0 = \begin{matrix} 0 \\ 1 \\ 2 \\ \vdots \\ r-2 \\ r-1 \end{matrix} \begin{pmatrix} I & & & & \\ p_j & \alpha p_j^{r-1} & & & \\ p_j^2 & & \alpha p_j^{r-2} & & \\ \vdots & & & \ddots & \\ p_j^{r-2} & & & & p_j^2 \\ p_j^{r-1} & & & & p_j \end{pmatrix}$$

where $\alpha \neq 0, 1$ is an element of the finite field and is multiplied to the diagonal in rows $1, \dots, \lfloor \frac{r}{2} \rfloor$. And define A_j^i by cyclicly

which is invertible by Lemma 6. Similarly, we can take out the i_2 -th column and row, and so on, and each sub matrix is again invertible. Thus, any matrix A is invertible and Construction 2 is MDS. ■

For example, one can easily check that the code in Figure 3 is able to recover the information from any two nodes. Therefore it is an MDS code.

IV. SUMMARY

In this paper, we presented constructions of MDS array codes that achieve the optimal rebuilding ratio $1/r$, where r is the number of redundancy nodes. The new codes are constructed based on our previous construction in [12] and improve the efficiency of the rebuilding access.

Now we mention a couple of open problems. For example, if there are $k - 1$ systematic nodes and r parity nodes, then our code has r^k rows. Namely, the code length is limited, are there codes that are longer given the number of rows? For example, when $r = 2$, we know an optimal rebuilding ratio construction with r^k rows and k systematic nodes:

$$A_j^0 = \begin{bmatrix} I & 0 \\ p_j & I \end{bmatrix}, A_j^1 = \begin{bmatrix} I & p_j \\ 0 & I \end{bmatrix}.$$

Here A_j^0, A_j^1 are the matrices that generate the parities, and we can take all $j \in [1, k]$. On the other hand, given r^k rows, it can be proven that any systematic and linear code with optimal ratio has no more than $k + 1$ systematic nodes. Thus the proposed code length can be improved by at most 2 nodes.

Finally, using the code in [12] one is able to rebuild any $e, 1 \leq e \leq r$, *systematic* erasures with an access ratio of e/r . However, it is an open problem to construct a code that can rebuild any e erasures with optimal access.

ACKNOWLEDGMENT

We thank Dimitris Papailiopoulos, Alexandros Dimakis and Viveck Cadambe for the inspiring discussions.

REFERENCES

- [1] N. Alon, "Combinatorial nullstellensatz," *Combinatorics Probability and Computing*, vol. 8, no. 1-2, pp. 7–29, Jan 1999.
- [2] M. Blaum, J. Brady, J. Bruck, and J. Menon, "EVENODD: an efficient scheme for tolerating double disk failures in RAID architectures," *IEEE Trans. on Computers*, vol. 44, no. 2, pp. 192–202, Feb. 1995.
- [3] V. R. Cadambe, C. Huang, S. A. Jafar, and J. Li, "Optimal repair of MDS codes in distributed storage via subspace interference alignment," Tech. Rep. arXiv:1106.1250, 2011.
- [4] V. Cadambe, S. Jafar, and H. Maleki, "Distributed data storage with minimum storage regenerating codes - exact and functional repair are asymptotically equally efficient," in *WINC*, 2010.
- [5] P. Corbett, B. English, A. Goel, T. Gnanac, S. Kleiman, J. Leong, and S. Sankar, "Row-diagonal parity for double disk failure correction," in *Proc. of the 3rd USENIX Symposium on File and Storage Technologies (FAST 04)*, 2004.
- [6] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems," *IEEE Trans. on Information Theory*, vol. 56, no. 9, pp. 4539–4551, 2010.
- [7] D. S. Papailiopoulos, A. G. Dimakis, and V. R. Cadambe, "Repair optimal erasure codes through hadamard designs," Tech. Rep. arXiv:1106.1634, 2011.
- [8] K. V. Rashmi, N. B. Shah, P. V. Kumar, "Optimal exact-regenerating codes for distributed storage at the MSR and MBR points via a product-matrix construction," Tech. Rep. arXiv:1005.4178, 2010.
- [9] N. B. Shah, K. V. Rashmi, P. V. Kumar, and K. Ramchandran, "Interference alignment in regenerating codes for distributed storage: necessity and code constructions," Tech. Rep. arXiv:1005.1634, 2010.
- [10] C. Suh and K. Ramchandran, "Exact regeneration codes for distributed storage repair using interference alignment," in *ISIT*, 2010.
- [11] C. Suh and K. Ramchandran, "On the existence of optimal exact-repair MDS codes for distributed storage," Tech. Rep. arXiv:1004.4663, 2010.
- [12] I. Tamo, Z. Wang, and J. Bruck, "MDS array codes with optimal rebuilding," Tech. Rep. arXiv:1103.3737, 2011.