

# Rack-Aware MSR Codes with Multiple Erasure Tolerance

Jiaojiao Wang, Dabin Zheng\*, Shenghua Li

## Abstract

The minimum storage rack-aware regenerating (MSRR) code is a variation of regenerating codes that achieves the optimal repair bandwidth for a single node failure in the rack-aware model. The authors in [1] and [25] provided explicit constructions of MSRR codes for all parameters to repair a single failed node. This paper generalizes the results in [1] to the case of multiple node failures. We propose a class of MDS array codes and scalar Reed-Solomon (RS) codes, and show that these codes have optimal repair bandwidth and error resilient capability for multiple node failures in the rack-aware storage model. Besides, our codes keep the same access level as the low-access constructions in [1] and [25].

**Keywords:** Distributed storage; Multiple erasure tolerance; MSRR codes; Universally error-resilient repair

## 1 Introduction

Maximum distance separable (MDS) codes are widely used in modern large-scale distributed storage system for which provide the maximum failure tolerance for a given amount of storage overhead. An important metric of repair efficiency is the *repair bandwidth*, i.e., the amount of data downloaded from other nodes for the purpose of the repair. Dimakis et al. in [3] have given a bound on the minimum number of symbols required for repair of a single failed node, which is called the cut-set bound of the repair bandwidth. A repair scheme that attains this bound is called optimal and such codes are said to be *minimum storage regenerating* (MSR) codes [3]. Over the last decade, important progresses have been made in the study of MSR codes, for example see [3, 4, 14, 16, 19, 20, 22, 23] and reference therein. In addition, the basic repair problem of MDS codes has been extended to the case that some of the helper nodes provide erroneous information [2, 12, 15, 22].

In a homogeneous distributed storage model, all nodes as well as communication between them are treated equally. However, modern data centers often have hierarchical topologies by organizing nodes in racks, where the cross-rack communication cost is much more expensive than the intra-rack communication cost. This motivates people to study the repair problem for hierarchical data centers and many progresses have been made recently [1, 5, 7–10, 13, 17, 18, 21, 24, 25]. In this paper, we focus on the constructions and repair schemes of MSR codes in the rack-aware storage model.

Let  $\mathcal{C}$  be an  $(n, k, \ell)$  array code over a finite field  $F$ , i.e., a collection of codewords  $c = (c_0, c_1, \dots, c_{n-1})$ , where each  $c_i$  is a vector of length  $\ell$  over  $F$ , or an element in the  $\ell$  extension of  $F$ . These code's coordinates are also called nodes. The amount of sub-packets stored in each node, i.e.,  $\ell$  is said to be *sub-packetization*. A code  $\mathcal{C}$  is called MDS if any  $k$  coordinates of the codeword suffice to recover its remaining  $n - k$  coordinates. Assume that  $k$  data blocks are encoded into a codeword of length  $n = u\bar{n}$  and stored across  $n$  nodes. The  $n$  nodes are organized equally into  $\bar{n}$  groups, also called racks, and every rack contains  $u$  nodes. The model of  $u = 1$  is the homogeneous case. To rule out the trivial case, we always assume that  $u \leq k$ , otherwise, a

---

\*Corresponding author. This work was partially supported by the National Natural Science Foundation of China under Grant Number 11971156.

Jiaojiao Wang, Dabin Zheng and Shenghua Li are with the Hubei Key Laboratory of Applied Mathematics, Faculty of Mathematics and Statistics, Hubei University, Wuhan 430062, China (E-mail: wjiaojiao@stu.hubu.edu.cn, dzheng@hubu.edu.cn, lish@hubu.edu.cn)

single node failure can be trivially repaired by  $u - 1$  surviving nodes within the same rack, and also assume that  $u \leq n - k$  to ensure that code has repair ability if an entire rack fails. The rack which contains the failed nodes is called the *host rack*.

The cut-set bound of repair bandwidth in a homogeneous distributed storage model has been generalized to the case of rack-aware storage model as follows [1, 8]. Let  $k = \bar{k}u + v$  ( $0 \leq v < u$ ), where  $u$  is the size of rack. Let  $\beta_u(h, \bar{d})$  denote the minimum number of symbols over  $F$  that one needs to download from the  $\bar{d}$ ,  $\bar{k} \leq \bar{d} \leq \bar{n} - 1$  helper racks to recover  $h$  failed nodes in the host rack. It was shown in [1] that

$$\beta_u(h, \bar{d}) \geq \frac{\bar{d}h\ell}{\bar{d} - \bar{k} + 1}. \quad (1)$$

A rack-aware repair scheme that achieves this bound is said to have optimal repair property and such codes are called *minimum storage rack-aware regenerating*(MSRR) codes for repairing  $h$  failed nodes in the rack-aware storage model.

In [22], Ye and Barg considered the error resilience capability in the repair process of multiple node failures in homogeneous distributed storage models. We generalize this concept to the case in rack-aware storage models. Let  $\bar{e}$  be a nonnegative integer with  $0 \leq \bar{e} \leq \lfloor \frac{\bar{d} - \bar{k}}{2} \rfloor$ . Suppose that a subset of  $\bar{e}$  racks out of  $\bar{d}$  helper racks provide erroneous information and define  $\beta_u(h, \bar{d}, \bar{e})$  to be the minimum number of symbols needed from the helper racks to repair the  $h$  failed nodes as long as the number of error racks in the helper racks is no more than  $\bar{e}$ . For  $\bar{k} + 2\bar{e} \leq \bar{d} \leq \bar{n} - 1$ , the bound in (1) can be generalized to the following,

$$\beta_u(h, \bar{d}, \bar{e}) \geq \frac{\bar{d}h\ell}{\bar{d} - 2\bar{e} - \bar{k} + 1}. \quad (2)$$

An  $(n, k, \ell)$  MDS code in rack-aware storage model is said to have *universally error-resilient* (UER)  $(h, \bar{d})$ -optimal repair property if the equality in (2) holds.

Recently, Hou et al. [8] studied the parameters and constructions of MSRR codes, but their constructions need some constraints on the parameters and the finite fields being large enough. The first explicit constructions of MSRR codes for all admissible parameters that support recovery of a single node failure were proposed by Chen et al. in [1]. Then Hou et al. [7] presented a coding framework that transformed an MSR code to an MSRR code. However, an MSRR code from such construction exists only if the finite field is sufficiently large. Zhou et al. [25] provided another class of MDS array codes for all parameters in the rack-aware model with smaller sub-packetization and size of underlying finite field. As far as we know, it remains an open problem to construct MSRR codes with error resilience capability for supporting recovery of multiple node failures. In this paper, we propose a class of MDS array codes and RS codes in the rack-aware storage model and corresponding repair schemes such that the codes have UER  $(h, \bar{d})$ -optimal repair property when the number of failed nodes  $h \leq u - v$ , where  $k = \bar{k}u + v$  ( $0 \leq v < u$ ). Moreover, we also provide repair schemes of discussed codes for the case  $h > u - v$ , and our schemes have the asymptotical UER  $(h, \bar{d} + 1)$ -optimal repair property. Comparisons between our MSRR array code and previous constructions are shown in Table 1.

Table 1: Comparisons with known MSRR codes, where  $\bar{s} = \bar{d} - 2\bar{e} - \bar{k} + 1$ ,  $k = \bar{k}u + v$ ,  $h$  is the largest error tolerance.

	sub-packet. $\ell$	$\bar{d}$	access per rack	$ F $	$h$	UER
[1]	$\bar{s}^{\bar{n}}$	$\bar{k} \leq \bar{d} \leq \bar{n} - 1$	$\frac{u\ell}{\bar{s}}$	$u( F  - 1),  F  \geq n + \bar{s} - 1$	1	-
[7]	$\bar{s}^{\lceil \frac{\bar{n}}{\bar{s}} \rceil}$	$\bar{d} = \bar{n} - 1$	$\frac{\ell}{\bar{s}} + (u - 1)\ell$	$ F  > k\ell \sum_{i=1}^{\min\{k, \bar{n}\}} \binom{n - \bar{n}}{k - i} \binom{\bar{n}}{i}$	1	-
[25]	$\bar{s}^{\lceil \frac{\bar{n}}{u - v} \rceil}$	$\bar{k} \leq \bar{d} \leq \bar{n} - 1$	$\frac{u\ell}{\bar{s}}$	$u( F  - 1),  F  > n$	1	-
this paper	$\bar{s}^{\bar{n}}$	$\bar{k} \leq \bar{d} \leq \bar{n} - 1$	$\frac{u\ell}{\bar{s}}$	$u( F  - 1),  F  > n$	$u - v$	✓

The rest of this paper is organized as follows. Section 2 proposes a class of MDS array codes and shows that they have the UER  $(h, \bar{d})$ -optimal repair property when the number of failed nodes  $h \leq u - v$ , and also

discusses the repair problem of the codes when  $h > u - v$  in the rack-aware storage system. In Section 3 we show that RS codes in [1] have the UER  $(h, \bar{d})$ -optimal repair property in the rack-aware storage system. Section 4 concludes this paper.

## 2 MSRR codes with multiple erasure tolerance

Let  $F$  be a finite field of size  $|F| > n$  and  $u$  denote the number of nodes in each rack satisfying  $u \mid (|F| - 1)$ . Let  $n = u\bar{n}$  and  $k = u\bar{k} + v$  for some  $v$  with  $0 \leq v < u$ . Let  $r = n - k$  and  $\bar{r} = \bar{n} - \bar{k}$ . Let  $\bar{d}$  denote the number of helper racks and  $\bar{e}$  be the largest acceptable number of erroneous racks in the  $\bar{d}$  helper racks satisfying  $(u, \bar{s}) = 1$ , where  $\bar{s} = \bar{d} - 2\bar{e} - \bar{k} + 1$ . Let  $\xi$  be a primitive element of  $F$  and  $\gamma$  be an element of  $F$  with multiplicative order  $u$ . Denote  $\lambda_{i,0} = \xi^i$  and  $\lambda_{i,j} = 1$  for  $i \in \{0, 1, \dots, \bar{n} - 1\}$  and  $j \in \{1, 2, \dots, \bar{s} - 1\}$ . In this section, we propose a class of MDS array codes which are slightly modified codes in [22] and show that the codes have the UER  $(h, \bar{d})$ -optimal repair property in rack-aware storage model for  $h \leq u - v$ . The objective array code is defined as follows.

$$\mathcal{C} = \left\{ (C_0, C_1, \dots, C_{n-1}) : \sum_{i=0}^{\bar{n}-1} \sum_{g=0}^{u-1} A_{i,g}^t C_{iu+g} = 0, t = 0, 1, \dots, r-1 \right\}, \quad (3)$$

here,  $C_j$  is a column vector of length  $\ell = \bar{s}\bar{n}$  over  $F$ , which denotes the  $j$ -th node and

$$A_{i,g} = \gamma^g A_i, \quad A_i = \sum_{j=0}^{\ell-1} \lambda_{i,j_i} e_j e_{j(i,j_i \oplus 1)}^T, \quad i = 0, 1, \dots, \bar{n} - 1, \quad g = 0, 1, \dots, u - 1, \quad (4)$$

where  $\{e_j : j = 0, 1, \dots, \ell - 1\}$  is the standard basis of  $F^\ell$  over  $F$ ,  $\oplus$  denotes addition modulo  $\bar{s}$ ,  $j_i$  is the  $i$ -th term of the base  $\bar{s}$  expansion of  $j = (j_{\bar{n}-1}, \dots, j_1, j_0)$  and  $j(i, b) = (j_{\bar{n}-1}, \dots, j_{i+1}, b, j_{i-1}, \dots, j_0)$ ,  $b = 0, 1, \dots, \bar{s} - 1$ .

By the similar discussion in Theorem VII.4 in [22] we have the following proposition.

**Proposition 2.1** *The array code  $\mathcal{C}$  given in (3) and (4) satisfies the MDS property.*

To show the optimal repair property of the code  $\mathcal{C}$ , we need the following result.

**Lemma 2.2** (*[22]*) *For two integers  $n, \ell > 0$ , let  $M_0, M_1, \dots, M_{n-1}$  be  $\ell$  order square matrices. For any  $i, j \in \{0, 1, \dots, n-1\}$ ,  $M_i M_j = M_j M_i$  and  $M_i - M_j$  is invertible, where  $i \neq j$ . Then*

$$M = \begin{pmatrix} I_\ell & \cdots & I_\ell \\ M_0 & \cdots & M_{n-1} \\ & \vdots & \\ M_0^{n-1} & \cdots & M_{n-1}^{n-1} \end{pmatrix}$$

*is invertible.*

The following theorem is our main result in this section.

**Theorem 2.3** *If the number  $h$  of failed nodes located in the same rack satisfies  $0 < h \leq u - v$ , then the array code  $\mathcal{C}$  defined in (3) and (4) has the UER  $(h, \bar{d})$ -optimal repair property.*

*Proof.* First we repair the linear combination of the nodes in the host rack from the remaining  $\bar{n} - 1$  surviving racks. Denote  $\Lambda_{i,j_i,0} = 1$  and  $\Lambda_{i,j_i,t} = \lambda_{i,j_i} \lambda_{i,j_i \oplus 1} \cdots \lambda_{i,j_i \oplus (t-1)}$  for  $t = 1, 2, \dots, r - 1$ . By direct calculations, we have

$$A_{i,g}^t = \left( \sum_{j=0}^{\ell-1} \lambda_{i,j_i} \gamma^g e_j e_{j(i,j_i \oplus 1)}^T \right)^t = \sum_{j=0}^{\ell-1} \Lambda_{i,j_i,t} \gamma^{gt} e_j e_{j(i,j_i \oplus t)}^T.$$

The  $(iu + g)$ -th node  $C_{iu+g}$  can be rewritten as  $\sum_{j=0}^{\ell-1} c_{iu+g,j} e_j$ . By combining these with (3), we get

$$\sum_{j=0}^{\ell-1} \sum_{i=0}^{\bar{n}-1} \sum_{g=0}^{u-1} \gamma^{gt} \Lambda_{i,j_i,t} c_{iu+g,j(i,j_i \oplus t)} e_j = 0, \quad t = 0, 1, \dots, r-1. \quad (5)$$

Without loss of generality, suppose that the  $(\bar{n} - 1)$ -th rack is the host rack. The equality (5) is rewritten coordinate wise as the following,

$$\Lambda_{\bar{n}-1, j_{\bar{n}-1}, t} \sum_{g=0}^{u-1} \gamma^{gt} c_{(\bar{n}-1)u+g, j(\bar{n}-1, j_{\bar{n}-1} \oplus t)} = - \sum_{i=0}^{\bar{n}-2} \Lambda_{i, j_i, t} \sum_{g=0}^{u-1} \gamma^{gt} c_{iu+g, j(i, j_i \oplus t)}, \quad (6)$$

where  $t = 0, 1, \dots, r-1$  and  $j = 0, 1, \dots, \ell-1$ . Consider the parity-check equations in (6) for  $t \in \{m, u + m, \dots, (\bar{s} - 1)u + m\}$ , where  $m$  is a fixed number in  $\{0, 1, \dots, u - v - 1\}$ . Since  $\gamma^u = 1$ , from (6) we have

$$\Lambda_{\bar{n}-1, j_{\bar{n}-1}, wu+m} \sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g, j(\bar{n}-1, j_{\bar{n}-1} \oplus (wu+m))} = - \sum_{i=0}^{\bar{n}-2} \Lambda_{i, j_i, wu+m} \sum_{g=0}^{u-1} \gamma^{gm} c_{iu+g, j(i, j_i \oplus (wu+m))}, \quad (7)$$

where  $w = 0, 1, \dots, \bar{s} - 1$  and  $j = 0, 1, \dots, \ell - 1$ . So,  $\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g, j(\bar{n}-1, j_{\bar{n}-1} \oplus (wu+m))}$  can be determined by  $\{\sum_{g=0}^{u-1} \gamma^{gm} c_{iu+g, j(i, j_i \oplus (wu+m))} : i = 0, 1, \dots, \bar{n} - 2\}$ . Since  $(u, \bar{s}) = 1$  and  $m < u$ , for any fixed  $j_i$  and  $m$ , we have that  $\{j_i \oplus (wu + m) : w = 0, 1, \dots, \bar{s} - 1\} = \{0, 1, \dots, \bar{s} - 1\}$ . Set  $\ell' = \bar{s}^{\bar{n}-1}$  and  $j_{\bar{n}-1} = 0$  in (7), then we have that  $\{\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g, j} : j = 0, 1, \dots, \ell - 1\}$  can be derived from  $\{\sum_{g=0}^{u-1} \gamma^{gm} c_{iu+g, j} : i = 0, 1, \dots, \bar{n} - 2; j = 0, 1, \dots, \ell' - 1\}$ , that is to say, the linear combination of nodes in the  $(\bar{n} - 1)$ -th rack can be determined by the vectors of length  $\ell'$  in the first  $\bar{n} - 1$  racks.

Then we prove that any  $\bar{d}$  ( $k \leq \bar{d} \leq \bar{n} - 1$ ) helper racks out of  $\bar{n} - 1$  surviving racks can recover the values  $\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g, j}$ ,  $j = 0, 1, \dots, \ell - 1$ . To this end, set  $C_{iu+g}^{(\ell')} = (c_{iu+g,0}, c_{iu+g,1}, \dots, c_{iu+g, \ell'-1})^T$ , where  $i = 0, 1, \dots, \bar{n} - 2$  and define

$$\mathcal{C}_m = \left( \sum_{g=0}^{u-1} \gamma^{gm} C_g^{(\ell')}, \sum_{g=0}^{u-1} \gamma^{gm} C_{u+g}^{(\ell')}, \dots, \sum_{g=0}^{u-1} \gamma^{gm} C_{(\bar{n}-2)u+g}^{(\ell')} \right), \quad m = 0, 1, \dots, u - v - 1. \quad (8)$$

Assume that there are  $\bar{d}$  helper racks and at most  $\bar{e}$  out of  $\bar{d}$  helper racks with  $\bar{k} + 2\bar{e} \leq \bar{d} \leq \bar{n} - 1$  have errors. Next, we show that  $\mathcal{C}_m$  is an  $(\bar{n} - 1, \bar{d} - 2\bar{e}, \ell')$  MDS array code for any fixed  $m \in \{0, 1, \dots, u - v - 1\}$ , that is to say, any  $\bar{d} - 2\bar{e}$  columns can represent all columns in  $\mathcal{C}_m$ .

Consider the parity-check equations of  $\mathcal{C}$  in (3) for  $t \in \{m, u + m, \dots, u(\bar{r} - \bar{s} - 1) + m\}$ , then

$$\sum_{i=0}^{\bar{n}-1} A_i^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{iu+g} = 0, \quad \eta = 0, 1, \dots, \bar{r} - \bar{s} - 1. \quad (9)$$

Since

$$A_i^{u\bar{s}} = \left( \sum_{j=0}^{\ell-1} \lambda_{i,j_i} e_j e_{j(i,j_i \oplus 1)}^T \right)^{u\bar{s}} = \sum_{j=0}^{\ell-1} (\Lambda_{i,j_i,\bar{s}})^u e_j e_j^T = \xi^{iu} I_\ell,$$

from the parity-check equations of  $\mathcal{C}$  in (3) for  $t \in \{u\bar{s} + m, u(\bar{s} + 1) + m, \dots, u(\bar{r} - 1) + m\}$ , we have

$$\sum_{i=0}^{\bar{n}-1} A_i^{u(\eta+\bar{s})+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{iu+g} = \sum_{i=0}^{\bar{n}-1} \xi^{iu} A_i^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{iu+g} = 0, \quad \eta = 0, 1, \dots, \bar{r} - \bar{s} - 1. \quad (10)$$

Multiplying  $\xi^{(\bar{n}-1)u}$  on the both sides of (9) and then subtracting (10) we get

$$\sum_{i=0}^{\bar{n}-2} (\xi^{(\bar{n}-1)u} - \xi^{iu}) A_i^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{iu+g} = 0. \quad (11)$$

Substituting  $A_i$  in (4) into (11), we have

$$\begin{aligned}
0 &= \sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) \left( \sum_{j=0}^{\ell-1} \lambda_{i,j_i} e_j e_{j(i,j_i \oplus 1)}^T \right)^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} \left( \sum_{j=0}^{\ell-1} c_{iu+g,j} e_j \right) \\
&= \sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) \sum_{g=0}^{u-1} \gamma^{gm} \left( \sum_{j=0}^{\ell-1} \Lambda_{i,j_i,u\eta+m} c_{iu+g,j(i,j_i \oplus (u\eta+m))} e_j \right).
\end{aligned} \tag{12}$$

The equality (12) is rewritten coordinate wise as the following:

$$\sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) \sum_{g=0}^{u-1} \gamma^{gm} \Lambda_{i,j_i,u\eta+m} c_{iu+g,j(i,j_i \oplus (u\eta+m))} = 0, \quad j = 0, 1, \dots, \ell-1. \tag{13}$$

Let

$$B_i = \sum_{j=0}^{\ell'-1} \lambda_{i,j_i} e_j^{(\ell')} (e_{j(i,j_i \oplus 1)}^{(\ell')})^T, \quad i = 0, 1, \dots, \bar{n}-2, \tag{14}$$

where  $\{e_j^{(\ell')} : j = 0, 1, \dots, \ell'-1\}$  is a set of standard basis column vectors in  $F^{\ell'}$ . It is known that  $B_i$  is the leading principle submatrix of  $A_i$  with order  $\ell'$ . By direct verifications, from (13) we have

$$\begin{aligned}
&\sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) B_i^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{iu+g}^{(\ell')} \\
&= \sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) \sum_{g=0}^{u-1} \gamma^{gm} \left( \sum_{j=0}^{\ell'-1} \lambda_{i,j_i} e_j^{(\ell')} (e_{j(i,j_i \oplus 1)}^{(\ell')})^T \right)^{u\eta+m} \sum_{j=0}^{\ell'-1} c_{iu+g,j} e_j^{(\ell')} \\
&= \sum_{i=0}^{\bar{n}-2} \left( \xi^{(\bar{n}-1)u} - \xi^{iu} \right) \sum_{g=0}^{u-1} \gamma^{gm} \sum_{j=0}^{\ell'-1} \Lambda_{i,j_i,u\eta+m} c_{iu+g,j(i,j_i \oplus (u\eta+m))} e_j^{(\ell')} \\
&= 0,
\end{aligned} \tag{15}$$

where  $\eta = 0, 1, \dots, \bar{r} - \bar{s} - 1$ .

Choose  $\bar{d} - 2\bar{e}$  columns in  $\mathcal{C}_m$  with index set  $H = \{p_0, p_1, \dots, p_{\bar{d}-2\bar{e}-1}\}$ . Let  $\{q_0, q_1, \dots, q_{\bar{r}-\bar{s}-1}\} = \{0, 1, \dots, \bar{n}-2\} \setminus H$ . Then (15) can be rewritten as follows:

$$\sum_{i=0}^{\bar{r}-\bar{s}-1} \left( \xi^{(\bar{n}-1)u} - \xi^{q_i u} \right) B_{q_i}^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{q_i u+g}^{(\ell')} = - \sum_{j=0}^{\bar{d}-2\bar{e}-1} \left( \xi^{(\bar{n}-1)u} - \xi^{p_j u} \right) B_{p_j}^{u\eta+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{p_j u+g}^{(\ell')}.$$

These  $\bar{r} - \bar{s}$  equations are rewritten as the matrix equation form as follows:

$$\begin{aligned}
&\underbrace{\begin{pmatrix} I_{\ell'} & \cdots & I_{\ell'} \\ B_{q_0}^u & \cdots & B_{q_{\bar{r}-\bar{s}-1}}^u \\ \vdots & & \vdots \\ B_{q_0}^{(\bar{r}-\bar{s}-1)u} & \cdots & B_{q_{\bar{r}-\bar{s}-1}}^{(\bar{r}-\bar{s}-1)u} \end{pmatrix}}_B \times \underbrace{\begin{pmatrix} (\xi^{(\bar{n}-1)u} - \xi^{q_0 u}) B_{q_0}^m & & \\ & \ddots & \\ & & (\xi^{(\bar{n}-1)u} - \xi^{q_{\bar{r}-\bar{s}-1} u}) B_{q_{\bar{r}-\bar{s}-1}}^m \end{pmatrix}}_D \\
&\times \begin{pmatrix} \sum_{g=0}^{u-1} \gamma^{gm} C_{q_0 u+g}^{(\ell')} \\ \vdots \\ \sum_{g=0}^{u-1} \gamma^{gm} C_{q_{\bar{r}-\bar{s}-1} u+g}^{(\ell')} \end{pmatrix} = - \begin{pmatrix} \sum_{j=0}^{\bar{d}-2\bar{e}-1} (\xi^{(\bar{n}-1)u} - \xi^{p_j u}) B_{p_j}^m \sum_{g=0}^{u-1} \gamma^{gm} C_{p_j u+g}^{(\ell')} \\ \vdots \\ \sum_{j=0}^{\bar{d}-2\bar{e}-1} (\xi^{(\bar{n}-1)u} - \xi^{p_j u}) B_{p_j}^{u(\bar{r}-\bar{s}-1)+m} \sum_{g=0}^{u-1} \gamma^{gm} C_{p_j u+g}^{(\ell')} \end{pmatrix}.
\end{aligned} \tag{16}$$

By the definition of  $B_i$  in (14), it is easy to verify that  $B_i$  and  $B_i^u - B_j^u$  are invertible, and  $B_i^u B_j^u = B_j^u B_i^u$  for  $j \neq i$ . So, the matrix  $B$  is invertible by Lemma 2.2, and then the matrix  $B \times D$  is invertible. Therefore, the columns in  $\mathcal{C}_m$  with index set  $H$  can represent all columns in  $\mathcal{C}_m$ , that is to say,  $\mathcal{C}_m$  is an  $(\bar{n} - 1, \bar{d} - 2\bar{e}, \ell')$  MDS array code.

Choosing any  $\bar{d}$  columns from  $\mathcal{C}_m$  also constitutes a  $(\bar{d}, \bar{d} - 2\bar{e}, \ell')$  MDS array code, which is a punctured code of  $\mathcal{C}_m$ . This code is viewed as a linear code over the  $\ell'$  extension of  $F$ . Then its minimum distance is  $2\bar{e} + 1$ , and so can correct  $\bar{e}$  errors. Therefore, by downloading any  $\bar{d}$  nodes in  $\mathcal{C}_m$ , we can recover the entire codeword as long as the number of erroneous nodes among the  $\bar{d}$  helper nodes is not greater than  $\bar{e}$ , and further recover the linear combination of nodes in the  $(\bar{n} - 1)$ -th rack  $\{\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g,j} : j = 0, \dots, \ell - 1\}$  from (7).

Finally, we recover the  $h$  specific failed nodes in the host rack. Assume that the index set of failed nodes in the host rack (i.e., the  $(\bar{n} - 1)$ -th rack) is  $\mathcal{F} = \{g_1, g_2, \dots, g_h\}$  and  $\mathcal{T} = \{0, 1, \dots, u - 1\} \setminus \mathcal{F}$ . We have that  $\Delta_m = \sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g,j}$  is known for  $j = 0, \dots, \ell - 1$  and  $m = 0, 1, \dots, u - v - 1$ . Taking  $m = 0, 1, \dots, h - 1$  we get

$$\begin{pmatrix} 1 & \cdots & 1 \\ \gamma^{g_1} & \cdots & \gamma^{g_h} \\ \vdots & & \vdots \\ \gamma^{g_1(h-1)} & \cdots & \gamma^{g_h(h-1)} \end{pmatrix} \begin{pmatrix} c_{(\bar{n}-1)u+g_1,j} \\ c_{(\bar{n}-1)u+g_2,j} \\ \vdots \\ c_{(\bar{n}-1)u+g_h,j} \end{pmatrix} = \begin{pmatrix} \Delta_0 - \sum_{g \in \mathcal{T}} c_{(\bar{n}-1)u+g,j} \\ \Delta_1 - \sum_{g \in \mathcal{T}} \gamma^g c_{(\bar{n}-1)u+g,j} \\ \vdots \\ \Delta_{h-1} - \sum_{g \in \mathcal{T}} \gamma^{g(h-1)} c_{(\bar{n}-1)u+g,j} \end{pmatrix}, \quad j = 0, 1, \dots, \ell - 1. \quad (17)$$

Since  $\gamma^{g_i} \neq \gamma^{g_{i'}}$  for all  $i, i' \in \{0, 1, \dots, u - 1\}$  with  $i \neq i'$ , from the linear system (17), we can recover the failed nodes in the host rack with index set  $\mathcal{F}$ .

To recover these nodes we have downloaded  $\bar{d}h\ell' = \frac{\bar{d}\ell h}{s}$  symbols. This value is exactly the lower bound in (2), and so the discussed code has the UER  $(h, \bar{d})$ -optimal repair property.  $\square$

**Remark 2.4** *In the repair process described above, to recover the  $h$ ,  $0 < h \leq u - v$  failed nodes of the code  $\mathcal{C}$  defined in (3) and (4), we have accessed the symbols in the set  $\{c_{iu+g} : i \in \mathcal{R}; g = 0, \dots, u - 1; j = 0, \dots, \ell' - 1\}$ , where  $\mathcal{R}$  is the index set of  $\bar{d}$  helper racks. Then the total number of the accessed symbols is  $\frac{\bar{d}u\ell}{s}$ . So, this code has the same low-access property as that of the codes constructed in [1] and [25].*

The discussion above shows that the MDS array code  $\mathcal{C}$  defined in (3) and (4) has optimal repair property when the number of failed nodes in the host rack is no more than  $u - v$ . When the number of failed nodes  $h$  is greater than  $u - v$ , the code  $\mathcal{C}$  has asymptotical UER  $(h, \bar{d} + 1)$ -optimal repair property, i.e., when the number of helper racks  $\bar{d} + 1$  is large enough, the ratio between the amount of download symbols and the optimal bound in (2) approaches 1.

**Theorem 2.5** *If the number  $h$  of failed nodes located in the same rack satisfies  $u - v < h \leq u$ , then the repair bandwidth of the array code  $\mathcal{C}$  defined in (3) and (4) is less than  $\frac{(\bar{d}+1)\ell h}{s}$  and the code  $\mathcal{C}$  has an error correction capability.*

*Proof.* The repair scheme consists of the following two main steps.

(1) For  $m = 0, 1, \dots, u - v - 1$ , by a similar calculation in Theorem 2.3 we can recover  $\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g,j}$ ,  $j = 0, 1, \dots, \ell - 1$  from any  $\bar{d}$  out of  $\bar{n} - 1$  columns in  $\mathcal{C}_m$  defined in (8) as long as the number of helper racks where errors occur is no more than  $\bar{e}$ .

(2) For  $m = u - v, u - v + 1, \dots, h - 1$ , the values of  $\eta$  in (9) and (10) could be  $0, 1, \dots, \bar{r} - \bar{s} - 2$ . We can get  $\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g,j}$ ,  $j = 0, 1, \dots, \ell - 1$  from any  $\bar{d} + 1$  out of  $\bar{n} - 1$  columns in  $\mathcal{C}_m$  as long as the number of helper racks where errors occur is no more than  $\bar{e}$ . Repeating the calculations in (11)-(16), we can show that  $\mathcal{C}_m$  defined in (8) for  $m \in \{u - v, u - v + 1, \dots, h - 1\}$  is also an  $(\bar{n} - 1, \bar{d} + 1 - 2\bar{e}, \ell')$  MDS array code. So, we can get  $\sum_{g=0}^{u-1} \gamma^{gm} c_{(\bar{n}-1)u+g,j}$ ,  $j = 0, 1, \dots, \ell - 1$  from any  $\bar{d} + 1$  out of  $\bar{n} - 1$  columns in  $\mathcal{C}_m$ . By a similar argument to that in Theorem 2.3, we know that this code can correct at most  $\bar{e}$  errors among  $\bar{d} + 1$  helper racks. Then, by the similar calculation in equation (17), we can recover the  $h$  failed nodes.

In the two steps above, we download  $\frac{\bar{d}\ell(u-v)}{\bar{s}}$  and  $\frac{(\bar{d}+1)\ell(h-u+v)}{\bar{s}}$  symbols from  $\bar{d}$  and  $\bar{d}+1$  helper racks, respectively. So, to recover the  $h$  failed nodes we have altogether downloaded  $\frac{\bar{d}\ell(u-v)}{\bar{s}} + \frac{(\bar{d}+1)\ell(h-u+v)}{\bar{s}} = \frac{\bar{d}\ell h}{\bar{s}} + \frac{\ell(h-u+v)}{\bar{s}}$  symbols. Note that  $v < u$ , then  $h-u+v < h$ . Therefore,  $\frac{\bar{d}\ell h}{\bar{s}} + \frac{\ell(h-u+v)}{\bar{s}} < \frac{(\bar{d}+1)\ell h}{\bar{s}}$ . In this case, the ratio of the amount of download symbols to the optimal repair bandwidth given in (2) is less than  $1 + \frac{1}{\bar{d}-2\bar{e}-\bar{k}+1}$ . So, the repair bandwidth of the code approaches optimal level if the number of helper racks is large enough when  $h > u - v$ .  $\square$

### 3 Rack-aware RS codes with multiple erasure tolerance

Reed-Solomon codes, the most practically used MDS codes, have been employed in many distributed storage systems. Guruswami and Wootters first proposed the optimal repair scheme of RS codes for the homogeneous distributed storage system [6]. Chen and Barg in [1] studied the repair procedure of RS codes having optimal repair property for a single node failure in the rack-aware storage system. This section generalizes the repair process of RS codes in [1] to the case of the multiple node failures. We propose a new repair scheme for the RS codes defined in [1] which can optimally repair multiple failed nodes from arbitrary  $\bar{d}$  helper racks. Moreover, the error correction capability and the low-access property of these RS codes are discussed.

Let  $q$  be a power of a prime and  $\mathbb{F}_q$  be a finite field of  $q$  elements. Let  $u$  be the size of the rack with  $u|(q-1)$  and  $n = \bar{n}u$ . Let  $k = \bar{k}u + v$ ,  $0 \leq v < u$ . Let  $\bar{d}$  denote the number of helper racks and  $\bar{e}$  ( $0 \leq \bar{e} \leq \lfloor \frac{\bar{d}-2\bar{k}}{2} \rfloor$ ) be the largest acceptable number of erroneous racks in the  $\bar{d}$  helper racks. Let  $\bar{s} = \bar{d} - 2\bar{e} - \bar{k} + 1$  and  $p_0, \dots, p_{\bar{n}-1}$  be  $\bar{n}$  distinct primes such that

$$p_i \equiv 1 \pmod{\bar{s}} \text{ and } p_i > u, \quad i = 0, 1, \dots, \bar{n} - 1.$$

Let  $\gamma$  be an element in  $\mathbb{F}_q$  with the multiplicative order  $u$  and  $\lambda_i$  an element of degree  $p_i$  over  $\mathbb{F}_q$  for  $i = 0, 1, \dots, \bar{n} - 1$ . Let

$$F_i := \mathbb{F}_q(\lambda_j : j \in \{0, 1, \dots, \bar{n} - 1\} \setminus \{i\}), \quad i = 0, 1, \dots, \bar{n} - 1 \text{ and } F := \mathbb{F}_q(\lambda_0, \lambda_1, \dots, \lambda_{\bar{n}-1}).$$

Let  $K$  be an extension of  $F$  of degree  $\bar{s}$  and  $\alpha \in K$  be a generating element of  $K$  over  $F$ , i.e.,  $K = F(\alpha)$ . Thus, for any  $i \in \{0, 1, \dots, \bar{n} - 1\}$  we have the chain of inclusions

$$\mathbb{F}_q \subset F_i \subset F \subset K.$$

So,  $K$  is the  $\ell$ -th degree extension of  $\mathbb{F}_q$ , where  $\ell = \bar{s} \prod_{i=0}^{\bar{n}-1} p_i$  is called sub-packetization. Choosing the set of evaluation points as  $\Omega = \{\lambda_0\gamma^0, \dots, \lambda_0\gamma^{u-1}, \dots, \lambda_{\bar{n}-1}\gamma^0, \dots, \lambda_{\bar{n}-1}\gamma^{u-1}\}$ , we define the objective RS code as follows:

$$C = \left\{ (f(\lambda_0\gamma^0), \dots, f(\lambda_0\gamma^{u-1}), \dots, f(\lambda_{\bar{n}-1}\gamma^0), \dots, f(\lambda_{\bar{n}-1}\gamma^{u-1})) : f(x) \in K[x], \deg(f) < k \right\}. \quad (18)$$

To describe the repair procedure for multiple node failures in rack-aware model, we first recall some necessary preliminaries.

**Lemma 3.1** ([20]) *For  $i \in \{0, 1, \dots, \bar{n} - 1\}$ , consider the  $F_i$ -linear subspace  $S_i$ ,*

$$S_i = \text{Span}_{F_i} \left\{ \sum_{m=0}^{\bar{s}-1} \alpha^m \lambda_i^{u(p_i-1)}, \alpha^j \lambda_i^{u(j+t\bar{s})}, \quad j = 0, 1, \dots, \bar{s} - 1, \quad t = 0, 1, \dots, \frac{p_i - 1}{\bar{s}} - 1 \right\}.$$

Then

$$\dim_{F_i} S_i = p_i, \quad S_i + S_i \lambda_i^u + \dots + S_i \lambda_i^{u(\bar{s}-1)} = K,$$

where  $S_i \beta = \{\theta \beta, \theta \in S_i\}$ , and the operation  $+$  is the Minkowski sum of sets,  $T_1 + T_2 = \{t_1 + t_2 : t_1 \in T_1, t_2 \in T_2\}$ .



**Lemma 3.2** ([11]) *Let  $\mathbb{F}$  be a finite field and  $\mathbb{K}$  an  $m$ -degree extension of  $\mathbb{F}$ . Let  $\{\xi_1, \xi_2, \dots, \xi_m\}$  be a basis of  $\mathbb{K}$  over  $\mathbb{F}$ , and  $\{\bar{\xi}_1, \bar{\xi}_2, \dots, \bar{\xi}_m\}$  its dual basis. Then for any  $\beta \in \mathbb{K}$ , we have*

$$\beta = \sum_{t=1}^m \text{tr}_{\mathbb{K}/\mathbb{F}}(\xi_t \beta) \bar{\xi}_t.$$

**Theorem 3.3** *If the number  $h$  of failed nodes located in the same rack satisfies  $0 < h \leq u - v$ , then the RS code  $\mathcal{C}$  defined in (18) has the UER  $(h, \bar{d})$ -optimal repair property. The repair procedure accesses  $\ell/\bar{s}$  symbols on each of the nodes in the  $\bar{d}$  helper racks, and the repair scheme is independent of the choice of the subset of  $\bar{d}$  helper racks.*

*Proof.* The coordinates of every codeword in  $\mathcal{C}$  are viewed as vectors over the field  $\mathbb{F}_q$ , and they also represent data on each node. The node size equals  $\ell = \bar{s} \prod_{i=0}^{\bar{n}-1} p_i$ , which is the extension degree of  $K$  over  $\mathbb{F}_q$ .

Let  $\mathcal{C}^\perp$  denote the dual code of  $\mathcal{C}$ . It is known that  $\mathcal{C}^\perp$  is a generalized RS code, which has the following form:

$$\mathcal{C}^\perp = \left\{ (v_{0,0}f(\lambda_0\gamma^0), \dots, v_{0,u-1}f(\lambda_0\gamma^{u-1}), \dots, v_{\bar{n}-1,0}f(\lambda_{\bar{n}-1}\gamma^0), \dots, v_{\bar{n}-1,u-1}f(\lambda_{\bar{n}-1}\gamma^{u-1})) : \right. \\ \left. f(x) \in K[x], \deg(f) < n - k \right\}, \quad (19)$$

where  $v_{i,j}$  are nonzero in  $K$  and can be represented by  $\lambda_i$  and  $\gamma^j$  for  $i \in \{0, 1, \dots, \bar{n}-1\}$  and  $j \in \{0, 1, \dots, u-1\}$ .

Denote  $r = n - k$ . Take  $f(x)$  in (19) to be  $x^t$ , where  $t = 0, 1, \dots, r - 1$ . For every codeword  $(C_0, C_1, \dots, C_{n-1}) \in \mathcal{C}$ , we have

$$\sum_{i=0}^{\bar{n}-1} \sum_{g=0}^{u-1} v_{i,g} (\lambda_i \gamma^g)^t C_{iu+g} = 0, \quad t = 0, 1, \dots, r - 1. \quad (20)$$

Assume that the  $i^*$ -th rack is the host rack. Let  $\{e_0, e_1, \dots, e_{p_{i^*}-1}\}$  be a basis of the vector space  $S_{i^*}$  over the field  $F_{i^*}$ . From (20) we have

$$e_j \sum_{i=0}^{\bar{n}-1} \sum_{g=0}^{u-1} v_{i,g} (\lambda_i \gamma^g)^t C_{iu+g} = 0, \quad t = 0, 1, \dots, r - 1, \quad (21)$$

where  $j \in \{0, 1, \dots, p_{i^*} - 1\}$ . Consider the subset of the parity-check equations in (21) with indices  $t = m, u + m, \dots, (\bar{r} - 1)u + m$  for some fixed  $m \in \{0, 1, \dots, u - v - 1\}$ , then

$$e_j \lambda_{i^*}^{uw+m} \sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g} = -e_j \sum_{i \neq i^*} \lambda_i^{uw+m} \sum_{g=0}^{u-1} v_{i,g} \gamma^{gm} C_{iu+g}, \quad w = 0, 1, \dots, \bar{r} - 1. \quad (22)$$

Let  $\text{tr}_{i^*}(\cdot) = \text{tr}_{K/F_{i^*}}(\cdot)$  be the trace mapping from  $K$  to  $F_{i^*}$ . Since  $\lambda_i \in F_{i^*}$  for  $i \neq i^*$  and  $\gamma \in \mathbb{F}_q$ , applying  $\text{tr}_{i^*}(\cdot)$  to the both sides of (22), we have

$$\text{tr}_{i^*} \left( e_j \lambda_{i^*}^{uw+m} \sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g} \right) = - \sum_{i \neq i^*} \lambda_i^{uw+m} \sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}), \quad w = 0, 1, \dots, \bar{r} - 1, \quad (23)$$

where  $j \in \{0, 1, \dots, p_{i^*} - 1\}$ . For a given  $m \in \{0, 1, \dots, u - v - 1\}$ , from (23) we can recover the set  $\{\text{tr}_{i^*}(e_j \lambda_{i^*}^{uw+m} \sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}) : w = 0, 1, \dots, \bar{r} - 1\}$  from  $\{\sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*}(e_j v_{i,g} C_{iu+g}) : i \in [0, \bar{n} - 1] \setminus \{i^*\}\}$ , where  $j \in \{0, 1, \dots, p_{i^*} - 1\}$ . Since  $\bar{s} = \bar{d} - 2\bar{e} - \bar{k} + 1 \leq \bar{r}$ , we can recover the set

$$\left\{ \text{tr}_{i^*}(e_j \lambda_{i^*}^{uw+m} \sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}) : w = 0, 1, \dots, \bar{s} - 1, j = 0, 1, \dots, p_{i^*} - 1 \right\} \quad (24)$$



from the set

$$\left\{ \sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}) : i \in \{0, 1, \dots, \bar{n} - 1\} \setminus \{i^*\}, j = 0, 1, \dots, p_{i^*} - 1 \right\}. \quad (25)$$

By Lemma 3.1, we know that  $\{e_j \lambda_{i^*}^{uw+m} : j = 0, 1, \dots, p_{i^*} - 1; w = 0, \dots, \bar{s} - 1\}$  is a basis of  $K$  over  $F_{i^*}$ . By choosing its dual basis, from Lemma 3.2 we can recover  $\sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}$  from the set in (24). So, obtaining the values in the set in (25) efficiently is the key to recovering the linear combination of the failed nodes in the host rack, i.e.,  $\sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}$ .

Next, we will discuss the process to repair the set in (25) with the minimum amount of download symbols in the case of errors in helper racks. For all  $i \in \{0, 1, \dots, \bar{n} - 1\} \setminus \{i^*\}$  and  $m \in \{0, 1, \dots, u - v - 1\}$ , we define  $u - v$  array codes as follows:

$$\mathcal{C}_m = (\Upsilon_{m,0}, \dots, \Upsilon_{m,i^*-1}, \Upsilon_{m,i^*+1}, \dots, \Upsilon_{m,\bar{n}-1}), \quad m = 0, 1, \dots, u - v - 1, \quad (26)$$

where  $\Upsilon_{m,i} = (\sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}), j = 0, 1, \dots, p_{i^*} - 1)^\perp$ . By slight modifying the proof of the case in homogeneous storage case in [2, Sec.III] we can show that  $\mathcal{C}_m$  in (26) is an  $(\bar{n} - 1, \bar{d} - 2\bar{e}, p_{i^*})$  MDS array code for any  $m \in \{0, 1, \dots, u - v - 1\}$ . So, for an integer  $\bar{d}$  with  $\bar{k} + 2\bar{e} \leq \bar{d} \leq \bar{n} - 1$ , any  $\bar{d}$  out of  $\bar{n} - 1$  columns in  $\mathcal{C}_m$  suffice to recover all columns of  $\mathcal{C}_m$  as long as the number of errors in the  $\bar{d}$  columns is not greater than  $\bar{e}$ . Hence, in such case,  $\sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}$  can be recovered from  $\bar{d}$  columns in  $\mathcal{C}_m$  for any  $m \in \{0, 1, \dots, u - v - 1\}$ . By similar calculations to that in Theorem 2.3, we can recover the  $h$  failed nodes in the host rack from known linear combinations  $\sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}$ ,  $m = 0, 1, \dots, u - v - 1$ . Moreover, to repair  $h$  failed nodes in the same rack, we have downloaded  $\frac{\bar{d} \ell h}{\bar{s}}$  symbols over  $\mathbb{F}_q$  from the helper racks, and this meets the cut-set bound in (2). So, the RS code in (18) has the optimal repair bandwidth and error correction capability when the number of helper racks where the error occurred is no more than  $\bar{e}$ .

During the repair process described above, the amount of access symbols is  $\bar{d} u \ell$ . In the following we further discuss how to reduce this amount. Let  $p^* = \prod_{i=0}^{\bar{n}-1} p_i / p_{i^*}$  and let  $\{\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{p^*-1}\}$  be a basis of  $F_{i^*}$  over  $\mathbb{F}_q$  and  $\{\varepsilon_0^*, \varepsilon_1^*, \dots, \varepsilon_{p^*-1}^*\}$  its dual basis. By Lemma 3.2,

$$\begin{aligned} \sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}) &= \sum_{t=0}^{p^*-1} \text{tr}_{F_{i^*}/\mathbb{F}_q} \left( \varepsilon_t \left( \sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}) \right) \right) \varepsilon_t^* \\ &= \sum_{g=0}^{u-1} \gamma^{gm} \sum_{t=0}^{p^*-1} \text{tr}_{K/\mathbb{F}_q} (\varepsilon_t e_j v_{i,g} C_{iu+g}) \varepsilon_t^*, \end{aligned} \quad (27)$$

where  $j = 0, 1, \dots, p_{i^*} - 1$ . It is easy to show that the vectors  $\varepsilon_t e_j, t = 0, 1, \dots, p^* - 1, j = 0, 1, \dots, p_{i^*} - 1$  are in  $K$  and linearly independent over  $\mathbb{F}_q$ . We expand these vectors to a basis of  $K$  over  $\mathbb{F}_q$ , and denote it by  $\{\beta_i, i = 0, 1, \dots, \ell - 1\}$ , and its dual basis is represented as  $\{\beta_i^*, i = 0, 1, \dots, \ell - 1\}$ . For any  $i \neq i^*$  and  $g \in \{0, 1, \dots, u - 1\}$ , the element  $v_{i,g} C_{iu+g}$  can be represented as follows:

$$v_{i,g} C_{iu+g} = \sum_{b=0}^{\ell-1} c_{iu+g,b} \beta_b^*, \quad (28)$$

where  $c_{iu+g,b} \in \mathbb{F}_q$ . Substituting (28) into (27) we have

$$\begin{aligned} \sum_{g=0}^{u-1} \gamma^{gm} \text{tr}_{i^*} (e_j v_{i,g} C_{iu+g}) &= \sum_{g=0}^{u-1} \gamma^{gm} \sum_{t=0}^{p^*-1} \text{tr}_{K/\mathbb{F}_q} \left( \varepsilon_t e_j \sum_{b=0}^{\ell-1} c_{iu+g,b} \beta_b^* \right) \varepsilon_t^* \\ &= \sum_{g=0}^{u-1} \gamma^{gm} \sum_{t=0}^{p^*-1} \sum_{b=0}^{\ell-1} \text{tr}_{K/\mathbb{F}_q} (\varepsilon_t e_j \beta_b^*) c_{iu+g,b} \varepsilon_t^*. \end{aligned} \quad (29)$$

The equality (29) shows that  $\sum_{g=0}^{u-1} v_{i^*,g} \gamma^{gm} C_{i^*u+g}$  can be repaired by accessing the symbols  $c_{iu+g,b}$  for which  $\text{tr}_{K/\mathbb{F}_q}(\varepsilon_t e_j \beta_b^*) \neq 0$  for  $t = 0, 1, \dots, p^* - 1, j = 0, 1, \dots, p_{i^*} - 1$  and  $b = 0, 1, \dots, \ell - 1$ . Recall that all  $\varepsilon_t e_j$  are in set  $\{\beta_i, i = 0, 1, \dots, \ell - 1\}$  and its dual basis is  $\{\beta_b^*, b = 0, 1, \dots, \ell - 1\}$ . So, the amount of access symbols is  $\bar{d} u p^* p_{i^*} = \bar{d} u \ell / \bar{s}$ . This amount of access symbols is the same as the low-access construction in [1] and [25].  $\square$

When the number of failed nodes in host rack is more than  $u - v$ , the discussed code has asymptotical UER  $(h, \bar{d} + 1)$ -optimal repair property. Combining the repair process in Theorem 2.5 and Theorem 3.3, we have the following result.

**Theorem 3.4** *If the number  $h$  of failed nodes located in the same rack satisfies  $u - v < h \leq u - 1$ , then the repair bandwidth of the RS code  $\mathcal{C}$  defined in (18) is less than  $\frac{(\bar{d}+1)\ell h}{\bar{s}}$ . Moreover, this code has the error correction capability and accesses  $\ell/\bar{s}$  symbols on each of the nodes in  $\bar{d} + 1$  helper racks in the process of repair.*

## 4 Concluding remark

In this paper we proposed a class of MDS array codes and RS codes in the rack-aware storage system, and showed that they have the UER  $(h, \bar{d})$ -optimal repair property when the number of failed nodes  $h \leq u - v$ . When  $u - v < h \leq u$ , the discussed codes have asymptotical UER  $(h, \bar{d} + 1)$ -optimal repair property. It is worthy of further designing MSRR codes and the corresponding repair scheme such that they have smaller sub-packetization and optimal access property.

## References

- [1] Z. Chen and A. Barg, "Explicit constructions of MSR codes for clustered distributed storage: the rack-aware storage model," *IEEE Trans. Inf. Theory*, vol. 66, no. 2, pp. 886-899, 2020.
- [2] Z. Chen, M. Ye and A. Barg, "Enabling optimal access and error correction for the repair of Reed-Solomon codes," *IEEE Trans. Inf. Theory*, vol. 66, no. 12, pp. 7439-7456, 2020.
- [3] A. G. Dimakis, P. B. Godfrey, Y. Wu, M. J. Wainwright and K. Ramchandran, "Network coding for distributed storage systems," *IEEE Trans. Inf. Theory*, vol. 56, no. 9, pp. 4539-4551, 2010.
- [4] S. Goparaju, A. Fazeli and A. Vardy, "Minimum storage regenerating codes for all parameters," *IEEE Trans. Inf. Theory*, vol. 63, no. 10, pp. 6318-6328, 2017.
- [5] S. Gupta and V. Lalitha, "Rack-aware cooperative regenerating codes," in *Proc. Int. Symp. Inf. Theory and Its Appl.*, pp. 264-268, 2020.
- [6] V. Guruswami and M. Wootters, "Repairing Reed-Solomon codes," *IEEE Trans. Inf. Theory*, vol. 63, no. 9, pp. 5684-5698, 2017.
- [7] H. Hou, P. Lee and Y. Han, "Minimum storage rack-aware regenerating codes with exact repair and small sub-packetization," *IEEE Int. Symp. Inf. Theory (ISIT)*, pp. 554-559, 2020.
- [8] H. Hou, P. Lee, K. Shum and Y. Hu, "Rack-aware regenerating codes for data centers," *IEEE Trans. Inf. Theory*, vol. 65, no. 8, pp. 4730-4745, 2019.
- [9] Y. Hu, P. Lee and X. Zhang, "Double regenerating codes for hierarchical data centers," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, pp. 245-249, 2016.
- [10] L. Jin, G. Luo and C. Xing, "Optimal repairing schemes for Reed-Solomon codes with alphabet sizes linear in lengths under the rack-aware model," arXiv:1911.08016[cs. IT], 2019.

- [11] R. Lidl and H. Niederreiter, *Finite Fields*. Cambridge, U.K.: Cambridge University. Press, 1984.
- [12] S. Pawar, S. El Rouayheb and K. Ramchandran, “Securing dynamic distributed storage systems against eavesdropping and adversarial attacks,” *IEEE Trans. Inf. Theory*, vol. 57, no. 10, pp. 6734-6753, 2011.
- [13] J. Pernas, C. Yuen, B. Gastón and J. Pujol, “Non-homogeneous two-rack model for distributed storage systems,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, pp. 1237-1241, 2013.
- [14] K. V. Rashmi, N. B. Shah and P. V. Kumar, “Optimal exact-regenerating codes for distributed storage at the MSR and MBR points via a product-matrix construction,” *IEEE Trans. Inf. Theory*, vol. 57, no. 8, pp. 5227-5239, 2011.
- [15] K. V. Rashmi, N. B. Shah, K. Ramchandran and P. Y. Kumar, “Regenerating codes for errors and erasures in distributed storage,” in *Proc. IEEE Int. Symp. Inf. Theory*, pp. 1202-1206, 2012.
- [16] N. Raviv, N. Silberstein and T. Etzion, “Constructions of high-rate minimum storage regenerating codes over small fields,” *IEEE Trans. Inf. Theory*, vol. 63, no. 4, pp. 2015-2038, 2017.
- [17] J.-Y. Sohn, B. Choi and J. Moon, “A class of MSR codes for clustered distributed storage,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, pp. 2366-2370, 2018.
- [18] J.-Y. Sohn, B. Choi, S. W. Yoon and J. Moon, “Capacity of clustered distributed storage,” *IEEE Trans. Inf. Theory*, vol. 65, no. 1, pp. 81-107, 2018.
- [19] I. Tamo, Z. Wang and J. Bruck, “Zigzag codes: MDS array codes with optimal rebuilding,” *IEEE Trans. Inf. Theory*, vol. 59, no. 3, pp. 1597-1616, 2012.
- [20] I. Tamo, M. Ye and A. Barg, “The repair problem for Reed-Solomon codes: optimal repair of single and multiple erasures with almost optimal node size,” *IEEE Trans. Inf. Theory*, vol. 65, no. 5, pp. 2673-2695, 2018.
- [21] M. A. Tebbi, T. H. Chan and C. W. Sung, “A code design framework for multi-rack distributed storage,” in *Proc. IEEE Inf. Theory Workshop (ITW)*, pp. 55-59, 2014.
- [22] M. Ye and A. Barg, “Explicit constructions of high-rate MDS array codes with optimal repair bandwidth,” *IEEE Trans. Inf. Theory*, vol. 63, no. 4, pp. 2001-2014, 2017.
- [23] M. Ye and A. Barg, “Explicit constructions of optimal-access MDS codes with nearly optimal sub-packetization,” *IEEE Trans. Inf. Theory*, vol. 63, no. 10, pp. 6307-6317, 2017.
- [24] Z. Zhang and L. Zhou, “Rack-Aware regenerating codes with multiple erasure tolerance,” arXiv:2106.03302[cs. IT], 2021.
- [25] L. Zhou and Z. Zhang, “Explicit construction of minimum storage rack-aware regenerating codes for all parameters,” arXiv:2103.15471 [cs. IT], 2021.