Minimum Storage Regenerating Codes For All Parameters

Sreechakra Goparaju

Arman Fazeli

Alexander Vardy

University of California San Diego, La Jolla, CA 92093, USA Email: {sgoparaju, afazelic, avardy}@ucsd.edu

Abstract—Regenerating codes for distributed storage have attracted much research interest in the past decade. Such codes trade the bandwidth needed to repair a failed node with the overall amount of data stored in the network. Minimum storage regenerating (MSR) codes are an important class of optimal regenerating codes that minimize (first) the amount of data stored per node and (then) the repair bandwidth. Specifically, an [n,k,d]-(α) MSR code $\mathbb C$ over $\mathbb F_q$ is defined as follows. Using such a code $\mathbb C$, a file $\mathcal F$ consisting of αk symbols over $\mathbb F_q$ can be distributed among n nodes, each storing α symbols, in such a way that:

- the file F can be recovered by downloading the content of any k of the n nodes; and
- the content of any failed node can be reconstructed by accessing any d of the remaining n-1 nodes and downloading $\alpha/(d-k+1)$ symbols from each of these nodes.

A common practical requirement for regenerating codes is to have the original file $\mathcal F$ available in uncoded form on some k of the n nodes, known as systematic nodes. In this case, several authors relax the defining node-repair condition above, requiring the optimal repair bandwidth of $d\alpha/(d-k+1)$ symbols for systematic nodes only. We shall call such codes systematic-repair MSR codes.

Unfortunately, explicit constructions of [n,k,d] MSR codes are known only for certain special cases: either low rate, namely $k/n \leqslant 0.5$, or high repair connectivity, namely d=n-1. Although setting d=n-1 minimizes the repair bandwidth, it may be impractical to connect to all the remaining nodes in order to repair a single failed node. Our main result in this paper is an explicit construction of systematic-repair [n,k,d] MSR codes for all possible values of parameters n,k,d. In particular, we construct systematic-repair MSR codes of high rate k/n>0.5 and low repair connectivity $k\leqslant d\leqslant n-1$. Such codes were not previously known to exist. In order to construct these codes, we solve simultaneously several repair scenarios, each of which is expressible as an interference alignment problem. Extension of our results beyond systematic repair remains an open problem.

I. Introduction

Distributed storage systems form the backbone for modern cloud computing, large—scale data servers, and peer—to—peer systems. The data in these systems is stored in a redundant fashion — typically via replication (for instance, Hadoop [1] and Google file systems [2] adopt a triple replication policy) — to safeguard data against not—so—infrequently occurring disk failures. An alternative approach to storing data on these systems, which highly reduces the redundancy involved in replication, is to use maximum distance separable (MDS) codes such as Reed—Solomon codes. Though MDS codes are the most

space–efficient for a targeted worst–case number of simultaneous node failures, they, unlike repetition codes, incur a high repair bandwidth¹ when the system undergoes the repair of a single node failure. A new class of erasure codes, called regenerating codes, was recently defined by Dimakis et al. [4] over a set of n nodes, which simultaneously optimizes storage efficiency, worst–case resilience and repair bandwidth for single node failures. These codes follow a trade–off curve which is intuitively evidenced by the contrast between repetition codes and MDS codes: the repair bandwidth decreases as the storage redundancy per node increases.

Formally, a file \mathcal{F} of size M, is said to be stored on a DSS consisting of n nodes, each with a storage capacity of α , using an [n,k,d]- (α) (or, in short, [n,k,d]) regenerating code, if it satisfies two properties:

- (a) *data recovery*: the file \mathcal{F} can be recovered using the contents of any k of the n nodes (this property will also be referred to as the *MDS property*); and
- (b) *repair property*: the contents of any node can be recovered using the contents of a *helper set* of any d other *helper* nodes, where each node transmits β number of symbols to the replacement node.

An *optimal* [n,k,d] regenerating code achieves the optimal value of total repair bandwidth $\gamma = d\beta$ (minimum repair bandwidth) for a given storage capacity α and M. This is given implicitly by the following trade–off:

$$M = \sum_{i=0}^{k-1} \min \left\{ \alpha, (d-i)\beta \right\}. \tag{1}$$

Most of the regenerating codes research (e.g. [5]–[14]) is focussed on the extremal points of this trade–off: MBR and MSR codes. *Minimum bandwidth regenerating* (MBR) codes achieve the optimal α when the repair bandwidth equals that of a repetition code. This paper concerns *minimum storage regenerating* (MSR) codes, often dubbed as *optimal bandwidth MDS codes*, because they are optimal regenerating codes that are *also* MDS codes². For these codes, $\alpha = M/k$, and the optimal repair bandwidth is given by:

$$\beta = \frac{\alpha}{d-k+1}.$$
 (2)

¹A recent work [3] revisits this for the case of Reed–Solomon codes, but we do not go into that here.

²To be precise, these are vector MDS codes, i.e., MDS codes over \mathbb{F}_q^{α} .

It is easy to see that the total repair bandwidth $d\beta$ is optimized when the number of helper nodes d=n-1. However, it is not always practical to connect to *all* the remaining nodes to aid the repair of a failed node. We therefore consider the following question: *Are there constructions of* [n,k,d] *MSR codes, for* d < n-1?

A. Previous Work

This question has not been wholly unanswered. The first MSR code constructions appeared in [6], [15], which roughly correspond to the family of parameters $\{n,k,d\}$ with rate $k/n \leqslant 1/2$. The asymptotic existence of MSR codes for all triples $\{n,k,d\}$ was eventually shown in [10] using interference alignment techniques developed for a wireless interference channel; these codes achieve optimality as a regenerating code (as well as approach the MSR point) *only* when $\alpha \to \infty$, i.e., $\beta/\alpha \to 1/(d-k+1)$, as $M \to \infty$.

MSR codes, being MDS vector codes, can be expressed as a set of k systematic vectors and n - k parity vectors (the corresponding nodes are referred to as systematic and parity nodes, respectively). For the high-rate $(k/n \ge 1/2)$ regime, code constructions were discovered independently in [11]-[13], [16] for the specific case of d = n - 1. Of these, the constructions in [11], [12], [16] focus on the relaxation of restricting optimal repair to systematic nodes in the system; we call the corresponding codes systematic-repair MSR codes. Practical systems usually store information in a systematic format. Parity nodes may fail, but as in the above works, we do not require optimal bandwidth repair for such nodes (maybe they are less urgent or critical). Clearly, any node can be repaired by reconstructing the whole file, so this covers the node repairability (even if suboptimally).

B. Contribution & Outline

We present the first³ high–rate finite– α constructions for systematic–repair MSR codes for d < n-1. We start by describing in Section II the representative code construction that contains the ideas behind those in [11], [12], [16]. Leveraging on this, we present our construction in Section III, but restrict to the case when the helper nodes contain the remaining k-1 systematic nodes. This restriction is removed in Section IV, thus rounding out the code construction. We conclude with some remarks in Section V.

II. Primer: Code Construction for d = n - 1

Let n = k + r denote the number of nodes in the distributed storage system, where each node has the capacity to store a vector of size α over \mathbb{F}_q . Throughout

³This work was first presented (invited) at the 53rd Annual Allerton Conference on Communication, Control, and Computing. A simultaneous result was presented at the same venue by Tamo and En Gad [17]. Recently and independently, Rawat et al. [18] have constructed MSR codes which optimally repair all nodes. However, the flavor of their construction, which is not systematic in nature, differs from ours.

this paper, we discuss systematic constructions and assume that the first k nodes are information nodes and store raw information, while the remaining r nodes correspond to the parities. We use the notation \mathbf{x}_i , $i \in [k]$, for the raw information vectors stored in the systematic nodes. The parity nodes are defined by

$$\mathbf{x}_{k+i} = \sum_{j=1}^{k} A_{ij} \mathbf{x}_{j}, \qquad i \in [r], \tag{3}$$

where A_{ij} 's are $\alpha \times \alpha$ encoding matrices. The generator matrix of the code is then given by

$$G = \begin{bmatrix} I & & 0 \\ & \ddots & \\ 0 & & I \\ A_{1,1} & \cdots & A_{1,k} \\ \vdots & \ddots & \vdots \\ A_{r,1} & \cdots & A_{r,k} \end{bmatrix} . \tag{4}$$

In this section, we consider MSR codes where d = n - 1. In other words, when a single node failure occurs, all the remaining nodes aid in its repair. We also restrict our attention to codes that consider failures only of the systematic nodes, and discuss in this section, a construction that underlies the ideas in [12], [16] and [11]. This construction will inform our generalization for the general parameter triple $\{n, k, d\}$ in Section III.

Remark: Wang et al. constructed an MSR code for d = n - 1 in [13] that achieves the optimal repair bandwidth also for parity nodes, albeit at the cost of some other metrics such as the number of symbols read from a node and the complexity of updating parities when systematic data changes. We leave for future the question of whether such a code exists when d < n - 1.

A commonly adopted strategy in constructing an MSR code is to first guarantee the optimal repair bandwidth property for a single failure (in this case, for a single systematic node failure), and then transform the construction to ensure the MDS property. This is illustrated in Example 1 below.

Example 1. Assume (n, k, d) = (4, 2, 3) and $\alpha = 4$. Let the first two nodes \mathbf{x}_1 and \mathbf{x}_2 be the systematic nodes, and let the parity nodes \mathbf{x}_3 and \mathbf{x}_4 be defined as

$$\mathbf{x}_{3} = \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{I} \mathbf{x}_{1} + \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{I} \mathbf{x}_{2},$$

$$\mathbf{x}_{4} = \underbrace{\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}}_{P_{1}} \mathbf{x}_{1} + \underbrace{\begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}}_{P_{2}} \mathbf{x}_{2}.$$

Figure 1.a depicts the component-wise storage in each node. It can be observed that a single failure in either x_1

C_1	C_2	C_3	C_4
$x_{1,1}$	$x_{2,1}$	$x_{1,1} + x_{2,1}$	$x_{1,3} + x_{2,2}$
$x_{1,2}$	$x_{2,2}$	$x_{1,2} + x_{2,2}$	$x_{1,4} + x_{2,1}$
$x_{1,3}$	$x_{2,3}$	$x_{1,3} + x_{2,3}$	$x_{1,1} + x_{2,4}$
$x_{1,4}$	$x_{2,4}$	$x_{1,4} + x_{2,4}$	$x_{1,2} + x_{2,3}$

C_1	C_2	C_3	C_4	
$x_{1,1}$	$x_{2,1}$	$x_{1,1} + x_{2,1}$	$x_{1,3} + 2x_{2,2}$	
$x_{1,2}$	$x_{2,2}$	$x_{1,2} + x_{2,2}$	$x_{1,4} + 2x_{2,1}$	
$x_{1,3}$	$x_{2,3}$	$x_{1,3} + x_{2,3}$	$x_{1,1} + 2x_{2,4}$	
$x_{1,4}$	$x_{2,4}$	$x_{1,4} + x_{2,4}$	$x_{1,2} + 2x_{2,3}$	
(b)				

Fig. 1. (a) Component wise storage in a (4,2,3) binary array code with optimal repair bandwidth for a single systematic node failure, described by $(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_1 + \mathbf{x}_2, P_1\mathbf{x}_1 + P_2\mathbf{x}_2)$; (b) A (4,2,3) MSR code in \mathbb{F}_5 described by $(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_1 + \mathbf{x}_2, P_1\mathbf{x}_1 + 2P_2\mathbf{x}_2)$. In both cases, gray cells are accessed to rebuild C_1 .

or \mathbf{x}_2 can be reconstructed by downloading $\alpha/2=2$ elements from each of the remaining d=3 nodes. However, the data is not recoverable if both \mathbf{x}_1 and \mathbf{x}_2 fail and hence, the code is not MDS. To overcome this problem, we associate a coefficient λ with P_2 such that $\begin{pmatrix} I & I \\ P_1 & \lambda P_2 \end{pmatrix}$ is non-singular. Note that,

$$\begin{vmatrix} I & I \\ P_1 & \lambda P_2 \end{vmatrix} = \det(\lambda P_2 - P_1) = \begin{vmatrix} 0 & \lambda & -1 & 0 \\ \lambda & 0 & 0 & -1 \\ -1 & 0 & 0 & \lambda \\ 0 & -1 & \lambda & 0 \end{vmatrix} = (\lambda^2 - 1)^2,$$

which is non-zero⁴ if q = 5 and $\lambda = 2$. Figure 1.b shows the component-wise storage for the resulting MSR code.

Construction 1 generalizes the construction given in Example 1 for an (n,k,n-1) MSR code. Note that any MSR code construction must specify both the generator matrix of the code as well as the optimal bandwidth repair strategy that is implemented on the code.

Construction 1. Let $\alpha = r^k$ and label the α elements $[0:r^k-1]$ by r-ary vectors in \mathbb{Z}_r^k . Define permutation f_j^ℓ on $[0:r^k-1]$ as follows:

$$f_j^{\ell}: \quad \mathbb{Z}_r^k \quad \to \quad \mathbb{Z}_r^k$$

$$\quad v \quad \mapsto \quad v + \ell e_j,$$

for $j \in [k]$ and $\ell \in [0:r-1] := \{0,1,\ldots,r-1\}$, where $\{e_1,e_2,\ldots,e_k\}$ is the standard vector basis for \mathbb{Z}_r^k . The mapping f_j^ℓ is bijective, and therefore, corresponds to a permutation on $[0:r^k-1]$. Let $P_{\ell,j}$ be the $\alpha \times \alpha$ matrix corresponding to the permutation f_j^ℓ , that is, $P_{\ell,j} \mathbf{x} = \mathbf{y}$, where $\mathbf{x}, \mathbf{y} \in \mathbb{F}_q^{\alpha}$, and $\mathbf{x}(v) = \mathbf{y}(f_j^{\ell}(v))$. In other words, $P_{\ell,j}$ scrambles the elements of a vector according to the permutation f_j^ℓ . (Notice that $P_{0,j} = I_{\alpha}$.)

1) $MSR\ Code$: The generator matrix of the code is given by (4), where $A_{i,j} = \lambda_{i,j} P_{i-1,j}$, $i \in [r]$ and $j \in [k]$. The non-zero coefficients $\lambda_{i,j} \in \mathbb{F}_q$ will be defined in Section II-B to ensure the MDS property.

⁴In general, if A,B,C, and D are nonsingular $\alpha \times \alpha$ matrices, then $\det \left(\begin{bmatrix} A & B \\ C & \lambda D \end{bmatrix} \right)$ is given by $\det(A) \det(\lambda D - CA^{-1}B)$, which is a polynomial of degree at most α in λ . If the field size is large enough, *i.e.* $q > \alpha$, one can always find a value for λ so that the 2×2 block matrix becomes non-singular as well. The same approach can be used to prove Lemma 5.

2) Repair Strategy: Let $Y_j = \{v \in [0, r^k - 1] : v \cdot e_j = 0\}$ denote a subset of $[0: r^k - 1]$. Y_j can be interpreted as those elements in $[0: \alpha - 1]$ whose label representation in \mathbb{Z}_r^k have a 0 in their j^{th} coordinate. If systematic node j fails, it is repaired by accessing the elements corresponding to Y_j from each of the remaining nodes, i.e., by accessing $\mathbf{x}_i(v)$, where $v \in Y_j$ and $j \neq i \in [n]$.

Construction 1 is obtained by first constructing an [n,k] array code⁵ (Section II-A) which guarantees the optimal bandwidth repair for a single systematic node failure. The array code is then transformed (Section II-B) to an MDS array code (and thereby, a systematic–repair MSR code) by transforming the encoding matrices of the parity nodes, while retaining the repair property.

A. Repair Property: Interference Alignment

The optimal repair bandwidth property of an [n, k, n-1] MSR code can be viewed as a *signal interference* problem: the objective is to retrieve the desired signal — the contents of the failed systematic node, say, \mathbf{x}_i — which, in the repair data downloaded from the remaining nodes, is interfered by partial contents of the remaining systematic nodes, \mathbf{x}_j , where $i \neq j \in [n]$. The solution, turns out to be an *interference alignment* strategy, where the repair data associated with the interfering systematic data is aligned, so as to minimize the interference. This is crystallized in the following lemma⁶.

Lemma 1. Let \mathbf{x}_i , $i \in [k]$, be the failed systematic node. For an [n,k,n-1] MSR code, the set of d=n-1 helper nodes is given by $\mathcal{D} = \{\mathbf{x}_j \mid j \in [n] \setminus \{i\}\}$. To recover the contents of the failed systematic node with the optimal repair bandwidth, it is necessary and sufficient to find n-1 (repair) matrices denoted by $\{S_j^i \in \mathbb{F}_q^{\alpha/r \times \alpha} \mid j \in [n] \setminus \{i\}\}$, where r=n-k, such that, for $j \in [k]$, $j \neq i$, the following two conditions are satisfied:

(a) signal recovery:

$$\operatorname{rank}\left(\begin{pmatrix} S_{k+1}^{i}A_{1,i} \\ S_{k+2}^{i}A_{2,i} \\ \vdots \\ S_{k+r}^{i}A_{r,i} \end{pmatrix}\right) = \alpha, \tag{5}$$

(b) interference alignment:

$$\operatorname{rank}\left(\begin{pmatrix} S_{j}^{i} \\ S_{k+1}^{i} A_{1,j} \\ \vdots \\ S_{k+r}^{i} A_{r,j} \end{pmatrix}\right) = \frac{\alpha}{r}. \tag{6}$$

Stated otherwise, to optimally repair \mathbf{x}_i , it is necessary and sufficient to find n-1 (repair) subspaces of dimension α/r ,

 $^5\mathrm{By}$ an [n,k] array code, we mean a set of k systematic vectors, and n-k parity vectors defined according to (3), which may or may not satisfy any properties.

⁶This result is known and has been used in several papers on MSR codes, but we state and prove it for completeness.

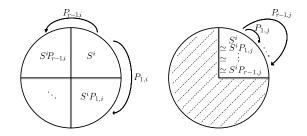


Fig. 2. Visualization of Lemma 1(a)[left], and Lemma 1(b)[right] to justify repair optimality in Construction 1.

denoted⁷ by $\{S_j^i | j \in [n] \setminus \{i\}\}$, where r = n - k, such that, for $j \in [k]$, $j \neq i$, the following two conditions are satisfied: (a) signal recovery:

$$S_{k+1}^i A_{1,i} \oplus \cdots \oplus S_{k+r}^i A_{r,i} \subseteq \mathbb{F}_a^{\alpha}, \tag{7}$$

(b) interference alignment:

$$S_j^i \simeq S_{k+s}^i A_{s,j}, \ \forall s \in [r],$$
 (8)

where \backsimeq denotes equality of subspaces, SA is the subspace obtained by operating the subspace S by the matrix A, and \oplus denotes the subspace sum.

For completeness, we provide a proof for Lemma 1 in Appendix A. Lemma 2 generalizes Lemma 1 when the number of helper nodes d < n - 1. This will be used later in Section III.

Lemma 2. (Corollary of Lemma 1.) In general, for an [n, k, d] MSR code, if the set of d = (k-1)+t < n-1 helper nodes is given by $\mathcal{D} = \{\mathbf{x}_j \mid j \in \mathcal{J} = [k] \setminus \{i\} \cup \{b_1, \cdots, b_t\}\}$ (where $b_i \in \{k+1, \ldots, n\}$ denote the t parity nodes in the helper set), it is necessary and sufficient to find d (repair) subspaces of dimension α/t denoted by $\{S_j^i \mid j \in \mathcal{J}\}$, such that, for $j \in [k]$, $j \neq i$, the following two conditions are satisfied:

$$S_{b_1}^i A_{b_1-k,i} \oplus \cdots \oplus S_{b_t}^i A_{b_t-k,i} \cong \mathbb{F}_q^{\alpha}, \tag{9}$$

$$S_j^i \simeq S_{b_s}^i A_{b_s-k,j}, \ \forall s \in [t].$$
 (10)

The optimal repair property of Construction 1 can now be justified.

Lemma 3. The repair strategy in Construction 1 is optimal with respect to repair bandwidth.

Proof: Define $S_j^i \triangleq S^i \triangleq Y_i, j \neq i$. Notice that the rank of subspace S^i is $r^{k-1} = \alpha/r$. Per definition, the permutation $P_{\ell,i}$ maps Y_i to $Y_i + \ell e_i = \{v \in [0, r^k - 1] : v \cdot e_i = \ell\}$. This implies that for any distinct $\ell, \ell' \in [0 : r - 1]$, the intersection $S^i P_{\ell,i} \cap S^i P_{\ell',i}$ contains only the all-zero vector. Thus the subspaces: $S^i, S^i P_{1,i}, \ldots, S^i P_{r-1,i}$, span

⁷Whenever this lemma is referenced, we use the subspace and matrix notation interchangeably as some proofs or expressions are clearer in one of the formats. We accordingly overload the notation S^i_j to refer to both the matrix and the subspace spanned by the row vectors of the matrix.

the space \mathbb{F}_q^{α} ($\alpha=r^k$) and the signal recovery condition(s) in Lemma 1 are satisfied. Furthermore, applying a permutation $P_{\ell,j}$ corresponding to a different coordinate $j\neq i$ maps Y_i to itself. This validates the interference alignment condition(s) in Lemma 1. Finally, note that the two conditions continue to be satisfied when replacing the permutations $P_{i-1,j}$ with any scaled versions $A_{i,j}=\lambda_{i,j}P_{i-1,j}$, because the scaling of the basis vectors does not change the relevant subspaces and thereby does not affect the conditions in Lemma 1.

B. MDS Property

This second step relies on the following two lemmas, the proofs of which are left to the reader.

Lemma 4. Let B denote the parity part of the generator matrix for an [n,k] array code denoted by \mathbb{C} , where

$$B = \left[\begin{array}{ccc} B_{1,1} & \cdots & B_{1,k} \\ \vdots & \ddots & \vdots \\ B_{r,1} & \cdots & B_{r,k} \end{array} \right].$$

Given that $B_{i,j}$ is non-singular for all i, j, then \mathbb{C} is an MDS array code if and only if any square sub-block-matrix B' of B is also non-singular, where

$$B' = \left[\begin{array}{ccc} B_{i_1,j_1} & \cdots & B_{i_1,j_t} \\ \vdots & \ddots & \vdots \\ B_{i_t,j_1} & \cdots & B_{i_t,j_t} \end{array} \right],$$

for some $\{i_1, \dots, i_t\} \subset [r], \{j_1, \dots, j_t\} \subset [k]$.

Lemma 5. Let B denote the $r\alpha \times k\alpha$ matrix associated with the parity part of the generator matrix for an [n,k] array code, as defined in Lemma 4. Given that $B_{i,j}$ is non-singular for all $i \in [r], j \in [k]$, and the field size q is large enough, there exist coefficients $\lambda_{i,j} \in \mathbb{F}_q$, such that all square sub-block-matrices of A are non-singular, where

$$A = \left[\begin{array}{ccc} \lambda_{1,1}B_{1,1} & \cdots & \lambda_{1,k}B_{1,k} \\ \vdots & \ddots & \vdots \\ \lambda_{r,1}B_{r,1} & \cdots & \lambda_{r,k}B_{r,k} \end{array} \right].$$

In other words, any parity generator matrix B for an [n,k] array code with non-singular encoding matrices can be transformed into a parity generator matrix A for an [n,k] MDS array code by multiplying the encoding matrices with appropriate scalar coefficients.

Proof Sketch: To obtain a valid set of $\lambda_{i,j}$'s, one may first sort the pairs (i,j) with respect to i+j increasingly, and then recursively choose a value for each $\lambda_{i,j}$ such that all sub-block-matrices with $\lambda_{i,j}A_{i,j}$ on their bottom right corner become non-singular. It suffices to have the field size q greater than the number of such sub-block-matrices at any step multiplied by α ;

$$|\mathbb{F}| > q_{\text{mds}} = \alpha \max_t \bigg\{ \binom{n-k-1}{t} \times \binom{k-1}{t} \bigg| t \in [k] \bigg\}.$$

III. CODE CONSTRUCTION FOR RESTRICTED HELPER SET

We now move to the construction of [n,k,d] systematic–repair MSR codes for any n,k, and d, where $k+1 \le d \le n-1$. In this section, we start with the restricted case when the helper set $\mathcal D$ includes all remaining k-1 systematic nodes. Let us begin with an example.

Example 2. Let us look at the case when [n,k,d] = [k+3,k,k+1] for $k \in \mathbb{N}$. Given a failure at the systematic node i, we are interested in repairing it by downloading $\frac{\alpha}{d-k+1} = \frac{\alpha}{2}$ symbols from each node in the helper set \mathcal{D}_i . Let us assume that \mathcal{D}_i includes all of the remaining k-1 systematic nodes. Hence, there are $\binom{3}{2} = 3$ different ways to choose \mathcal{D}_i depending on which two parity nodes are included in it. Let us use an indicator $a \in [3]$ to differentiate between these scenarios, and denote the helper set for each scenario by $\mathcal{D}_{i,a}$.

Construction. Let $\alpha = 2^{3k}$ and label the α elements $[0:2^{3k}-1]$ by binary vectors in \mathbb{Z}_2^{3k} . Define permutation f_j^{ℓ} on $[0:2^{3k}-1]$ as follows:

$$f_j^{\ell}: \quad \mathbb{Z}_2^{3k} \quad \rightarrow \quad \mathbb{Z}_2^{3k}$$
 $v \quad \mapsto \quad v + \ell e_{ij}$

for $j \in [3k]$ and $\ell \in \{0,1\}$, where $\{e_1,e_2,\cdots,e_{3k}\}$ is the standard vector basis for \mathbb{Z}_2^{3k} . The mapping f_j^ℓ is again bijective and therefore corresponds to a permutation on $[0:2^{3k}-1]$. As before, let $P_{\ell,j}$ be the $\alpha \times \alpha$ matrix corresponding to the permutation f_j^ℓ , that is, $P_{\ell,j}\mathbf{x} = \mathbf{y}$, where $\mathbf{x},\mathbf{y} \in \mathbb{F}_q^\alpha$, and $\mathbf{x}(v) = \mathbf{y}(f_j^\ell(v))$. (Notice again that $P_{0,j} = I_\alpha$.)

1) *MSR Code*: The generator matrix of the code is given by

$$G = \left[egin{array}{cccc} I & & & 0 \ & & \ddots & & \ 0 & & I \ A_{1,1} & \cdots & A_{1,k} \ A_{2,1} & \cdots & A_{2,k} \ A_{3,1} & \cdots & A_{3,k} \end{array}
ight],$$

where

$$A_{1,j} = \lambda_{1,j} \times P_{0,3j-2} \times P_{0,3j-1} \times I_{\alpha},$$

 $A_{2,j} = \lambda_{2,j} \times P_{1,3j-2} \times I_{\alpha} \times P_{0,3j},$
 $A_{3,j} = \lambda_{3,j} \times I_{\alpha} \times P_{1,3j-1} \times P_{1,3j},$ (11)

for $j \in [k]$. The non-zero coefficients $\lambda_{i,j} \in \mathbb{F}_q$ are again selected according to the discussion in Section II-B to establish the MDS property.

2) Repair Strategy via $\mathcal{D}_{i,1} = \{\mathbf{x}_j | j \in [k+3], j \neq i, k+3\}$: Let $Y_{i,1} = \{v \in [0, 2^{3k} - 1] : v \cdot e_{3i-2} = 0\}$ denote a subset of $[0: 2^{3k} - 1]$. $Y_{i,1}$ can be interpreted as those elements in $[0: 2^{3k} - 1]$ whose label representation in \mathbb{Z}_2^{3k} have a 0 in their $(3i-2)^{th}$ coordinate. If systematic node i fails, it can be repaired by accessing the elements corresponding to $Y_{i,1}$ from each of the helper nodes, i.e., by accessing $\mathbf{x}_j(v)$, where $v \in Y_{i,1}$ and $j \in \{1, 2, \cdots, i-1, i+1, \cdots, k, k+1, k+2\}$.

- 3) Repair Strategy via $\mathcal{D}_{i,2} = \{\mathbf{x}_j | j \in [k+3], j \neq i, k+2\}$: Similarly, let $Y_{i,2} = \{v \in [0, 2^{3k} 1] : v \cdot e_{3i-1} = 0\}$. If systematic node i fails, it can be repaired by accessing $\mathbf{x}_j(v)$, where $v \in Y_{i,2}$ and $j \in \{1, 2, \dots, i-1, i+1, \dots, k, k+1, k+3\}$.
- 4) Repair Strategy via $\mathcal{D}_{i,3} = \{\mathbf{x}_j | j \in [k+3], j \neq i, k+1\}$: Finally, let $Y_{i,3} = \{v \in [0,2^{3k}-1] : v \cdot e_{3i} = 0\}$ denote the location of the elements that have to get accessed if the systematic node i fails, i.e., node i can be repaired by accessing $\mathbf{x}_j(v)$, where $v \in Y_{i,3}$ and $j \in \{1,2,\cdots,i-1,i+1,\cdots,k,k+2,k+3\}$.

Justification of the repair strategy: Let \mathbf{x}_i , $i \in [k]$, be the failed systematic node. Define $Q_{u,v} = A_{u,v}\lambda_{u,v}^{-1}$, $u \in [3], v \in [k]$, which is a product of multiple permutation matrices, and hence can be viewed as a permutation matrix itself. In order to justify the repair strategy, it suffices to define the proper subspaces S_j^i that fulfill the two interference alignment conditions in Lemma 2. Let $U_{i,a}$ be the complimentary subset of $Y_{i,a}$ in \mathbb{Z}_2^{3k} , i.e.,

$$\begin{aligned} &U_{i,1} = \{v \in [0, 2^{3k} - 1] : v \cdot e_{3i-2} = 1\}, \\ &U_{i,2} = \{v \in [0, 2^{3k} - 1] : v \cdot e_{3i-1} = 1\}, \\ &U_{i,3} = \{v \in [0, 2^{3k} - 1] : v \cdot e_{3i} = 1\}. \end{aligned}$$

Given the code construction in (11), we can verify that

Now we define subspaces $S_{j,a}^i \triangleq S_a^i \triangleq Y_{i,a}, j \neq i, a \in [3]$. Let us for simplicity assume a = 1. The other scenarios follow the proof similarly. Based on (12), we observe that the permutation $Q_{2,i}$ maps the basis $Y_{i,1}$ to its complementary subset $U_{i,1}$ and vice versa, while $Q_{1,i}$ preserves both of them. Hence,

$$\operatorname{rank}\left(\left(\begin{array}{c}S_1^iQ_{1,i}\\S_1^iQ_{2,i}\end{array}\right)\right)=\operatorname{rank}\left(\left(\begin{array}{c}Y_{i,1}\\U_{i,1}\end{array}\right)\right)=\alpha.$$

Furthermore, $Y_{i,1}$ remains unchanged under any other permutation $Q_{t,j}$, $j \neq i$, and hence

$$\operatorname{rank}\left(\left(\begin{array}{c} S_1^i \\ S_1^i Q_{1,i'} \\ S_1^i Q_{2,i'} \end{array}\right)\right) = \operatorname{rank}\left(\left(\begin{array}{c} Y_{i,1} \\ Y_{i,1} \\ Y_{i,1} \end{array}\right)\right) = \frac{\alpha}{2}.$$

The key element in the construction is to satisfy the two requirements in Lemma 2 for any systematic failure and any such helper set \mathcal{D} . Let $\rho = d - k + 1$ denote the number of parity nodes in the helper set of size d. There are

$$\Omega = \begin{bmatrix} \omega_1 & \omega_2 & \omega_3 \\ \downarrow & \downarrow & \downarrow \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \Rightarrow Q_{1,j} = P_{0,3j-2} \times P_{0,3j-1} \times I_{\alpha}$$

$$Q_{2,j} = P_{1,3j-2} \times I_{\alpha} \times P_{0,3j}$$

$$Q_{3,j} = I_{\alpha} \times P_{1,3j-1} \times P_{1,3j}$$

Fig. 3. Relation between ω_a and (11).

 $\binom{r}{\rho}$ different ways to choose ρ parity nodes during the repair. Let us label these cases with numbers $a \in [\binom{r}{\rho}]$, and set \mathcal{R}_a to be the subset of parity nodes corresponding to case a.

Assume that $\mathcal{R}_a = \{\mathbf{x}_{k+d_1^{(a)}}, \mathbf{x}_{k+d_2^{(a)}}, \cdots, \mathbf{x}_{k+d_\rho^{(a)}}\}$ is the ordered representations, where $\{d_1^{(a)}, \cdots, d_\rho^{(a)}\} \subset [r]$. Finally, define r-ary vectors ω_a for $a \in [\binom{r}{\rho}]$ as

$$\omega_a(i) = \begin{cases} t-1 & \text{if } \exists t : i = d_t^{(a)}, \\ 0 & \text{otherwise.} \end{cases}$$

Construction 2. Let $\alpha = \rho^{k\binom{r}{\rho}}$ and label the α elements $[0:\alpha-1]$ by ρ -ary vectors in $\mathbb{Z}_{\rho}^{k\binom{r}{\rho}}$. Define permutation f_i^{ℓ} on $[0:\alpha-1]$ as follows:

$$\begin{array}{cccc} f_j^{\ell}: & \mathbb{Z}_{\rho}^{k\binom{r}{\rho}} & \to & \mathbb{Z}_{\rho}^{k\binom{r}{\rho}} \\ & v & \mapsto & v + \ell e_j, \end{array}$$

for $j \in [k\binom{r}{\rho}]$ and $\ell \in [0:\rho-1]$, where $\{e_1, \cdots, e_{k\binom{r}{\rho}}\}$ is the standard vector basis of $\mathbb{Z}_{\rho}^{k\binom{r}{\rho}}$. Let $P_{\ell,j}$ be the $\alpha \times \alpha$ matrix corresponding to the permutation f_j^{ℓ} .

1) *MSR Code*: The generator matrix of the [n, k, d] code is given by (4), where

$$A_{i,j} = \lambda_{i,j} \prod_{a \in [\binom{r}{\rho}]} P_{w_a(i), a + (j-1)\binom{r}{\rho}}, \text{ for } j \in [k], i \in [r].$$

The non-zero coefficients $\lambda_{i,j} \in \mathbb{F}_q$ are defined according to Section II-B to ensure the MDS property; and later will be modified again in Section IV.

2) Repair Strategy: Let \mathcal{R}_a correspond to the parity subset of the helper set \mathcal{D} . Define $Y_{j,a} \subset [0:\alpha-1]$ as $\{x \in [0,\alpha-1]: x \cdot e_{a+(j-1)\binom{r}{\rho}} = 0\}$. If systematic node j fails, it is repaired by accessing the elements corresponding to $Y_{j,a}$ from helper nodes, i.e., by accessing $\mathbf{x}_i(v)$, where $i \in \mathcal{D}$, and $v \in Y_{j,a}$.

Lemma 6. The repair strategy in Construction 2 is optimal with respect to repair bandwidth.

Proof: Let us first explain the role of ω_a by revisiting Example 2 via Figure 3. Here we assumed that

$$a = 1 \rightarrow \mathcal{R}_1 = \{\mathbf{x}_{k+1}, \mathbf{x}_{k+2}\} \rightarrow \omega_1 = (0, 1, 0)^t,$$

 $a = 2 \rightarrow \mathcal{R}_2 = \{\mathbf{x}_{k+1}, \mathbf{x}_{k+3}\} \rightarrow \omega_2 = (0, 0, 1)^t,$
 $a = 3 \rightarrow \mathcal{R}_3 = \{\mathbf{x}_{k+2}, \mathbf{x}_{k+3}\} \rightarrow \omega_3 = (0, 0, 1)^t.$

In general, the matrix $\Omega = \{\omega_1 | \omega_2 | \cdots | \omega_{\binom{r}{\rho}}\}$ is designed in a way that for any choice of $a \in [\binom{r}{\rho}]$ we can always find a column in Ω , denoted by ω_a , such that its intersection with r' rows associated with scenario a, forms $\{0,1,\cdots,\rho-1\}$.

Now assume that node i is failed and we are to perform an optimal systematic repair given parity repairs in $\mathcal{R}_a = \{\mathbf{x}_{k+d_1^{(a)}}, \mathbf{x}_{k+d_2^{(a)}}, \cdots, \mathbf{x}_{k+d_\rho^{(a)}}\}$. It is now clear that if we selected our subspaces as $S_{j,a}^i \triangleq S_a^i \triangleq Y_{i,a} = \{x | x \cdot e_{a+(i-1)\binom{r}{a}} = 0\}$, then

$$\begin{split} Y_{i,a}Q_{d_{1}^{(a)},i} &= & \{x|x\cdot e_{a+(i-1)\binom{r}{\rho}} = 0\},\\ Y_{i,a}Q_{d_{2}^{(a)},i} &= & \{x|x\cdot e_{a+(i-1)\binom{r}{\rho}} = 1\},\\ &\vdots\\ Y_{i,a}Q_{d_{\rho}^{(a)},i} &= & \{x|x\cdot e_{a+(i-1)\binom{r}{\rho}} = \rho - 1\}, \end{split}$$

and hence,

$$\operatorname{rank}\left(\left(\begin{array}{c}S_a^iQ_{d_1^{(a)},i}\\S_a^iQ_{d_2^{(a)},i}\\\vdots\\S_a^iQ_{d_n^{(a)},i}\end{array}\right)\right)=\rho\times\operatorname{rank}\left(Y_{i,a}\right)=\rho\frac{\alpha}{\rho}=\alpha.$$

The second condition in Lemma 2 is also automatically satisfied since

$$Y_{i,a} \simeq Y_{i,a} Q_{1,i'} \simeq Y_{i,a} Q_{2,i'} \simeq \cdots \simeq Y_{i,a} Q_{\binom{r}{\rho},i'} \quad \text{for } i' \neq i.$$

Lastly, we note that optimizing the sub-packetization parameter, α , is not the main concern. Although Construction 2 suggests a fairly large value, i.e. $\alpha = \rho^{k\binom{r}{\rho}}$, but it is clear that we do not need $\binom{r}{\rho}$ many columns in Ω to cover all the $\binom{r}{\rho}$ helper set selection scenarios. Indeed, α in Example 2 can be reduced to 2^{2k} , where $\Omega_{\text{new}} = \{\omega_1 | \omega_2\}$. We leave the optimizations of this kind to future work.

IV. CODE CONSTRUCTION FOR ANY HELPER SET

In this section, we show that Construction 2 in fact holds, even when an arbitrary set of d helper nodes is allowed to be chosen from the (n-1) surviving nodes. This generality merely imposes some additional constraints on the selection of the scaling coefficients $\lambda_{i,j}$ of the encoding matrices $A_{i,j} = \lambda_{i,j} Q_{i,j}$, where $Q_{i,j}$ is the (product) permutation matrix corresponding to $A_{i,j}$, as defined in Construction 2. We now arrive at the main theorem.

Theorem 7. Construction 2 gives an [n,k,d] systematic-repair MSR code for any set of d helper nodes, for a large enough field size for the scaling coefficients $\lambda_{i,j}$ for the encoding matrices $A_{i,j}$.

Proof: Part 1: First, we illustrate the proof by fixing d = k + 1, and taking an example set of helper nodes for an example failure of node \mathbf{x}_1 (or node 1). Let us denote the (indices of the) helper set by \mathcal{D} , and let $\mathcal{D} = \{h, h + 1, \ldots, k, k + 1, \ldots, k + h\}$, that is, there are h parity nodes and d - h = k + 1 - h systematic nodes in the helper set. Let $S_i^j(\mathcal{D})\mathbf{x}_i$ denote the repair information that node i sends to help in the repair of node j when \mathcal{D} is the set of helper nodes. (Wherever clear, we ignore the \mathcal{D} in the notation and simply write S_i^j .) When node 1 fails, the information we therefore have at its replacement node can be written as:

$$\begin{pmatrix} x_1 & x_2 & \cdots & x_h & x_{h+1} & \cdots & x_{k-1} & x_k \\ & & S_h^1 & & & & & \\ & & & S_{h+1}^1 & & & & \\ & & & & \ddots & & & \\ & & & & & S_{h+1}^1 & & \\ & & & & & & S_{h+1}^1 & & \\ S_{k+1}^1 A_{11} & S_{k+1}^1 A_{12} & \cdots & \cdots & S_{k+1}^1 A_{1k} & & \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots & \\ S_{k+h}^1 A_{h1} & S_{k+h}^1 A_{h2} & \cdots & \cdots & S_{k+h}^1 A_{hk} & & \\ \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_h \\ x_{h+1} \\ \vdots \\ x_{k-1} \\ x_k \end{pmatrix} (13)$$

Suppose all S_j^1 's in (13) be replaced by a repair subspace S^1 (corresponding to Lemma 1(b)) that we would have used if $\mathcal{D} = \{2,3,\ldots,k,k+1,k+2\}$. Specifically, suppose $S^1A_{1,1}$ and S^1 complete the space \mathbb{F}_q^n . Since S^1 and $S^1A_{i,j}$ denote the same subspace, for $j \neq 1$, the components of $\mathbf{x}_i, i \in \{h,h+1,\ldots,k\}$ can be easily subtracted from the information coming from the parity nodes , using that coming from the systematic nodes h to k. Thus, in order to recover \mathbf{x}_1 , we can concentrate on the following information at the replacement node:

$$\begin{pmatrix} S^{1}A_{11} & S^{1}A_{12} & \cdots & S^{1}A_{1,h-2} & S^{1}A_{1,h-1} \\ S^{1}A_{21} & S^{1}A_{22} & \cdots & S^{1}A_{2,h-2} & S^{1}A_{2,h-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ S^{1}A_{h1} & S^{1}A_{h2} & \cdots & S^{1}A_{h,h-2} & S^{1}A_{h,h-1} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{1} \\ \mathbf{x}_{2} \\ \vdots \\ \mathbf{x}_{h-2} \\ \mathbf{x}_{h-1} \end{pmatrix} \mathbf{1}_{4})$$

Let $S^1A_{i,j} = \lambda_{i,j}S^1Q_{i,j} = \lambda_{i,j}\widetilde{Q}_{i,j}S^1$, where $\widetilde{Q}_{i,j}$ is an $\alpha/2 \times \alpha/2$ matrix, and $(i,j) \neq (1,1)$. It must be noted that not only is S^1 dependent on the choice of \mathcal{D} , but so in turn is $\widetilde{Q}_{i,j}$. Let us also denote $S^1\mathbf{x}_i$ by $\widetilde{\mathbf{x}}_i$. Then, (14) can be rewritten as:

$$\begin{pmatrix} \lambda_{1,1}S^{1}Q_{1,1} & \lambda_{1,2}\widetilde{Q}_{1,2} & \cdots & \lambda_{1,h-1}\widetilde{Q}_{1,h-1} \\ \lambda_{2,1}\widetilde{Q}_{2,1}S^{1} & \lambda_{2,2}\widetilde{Q}_{2,2} & \cdots & \lambda_{2,h-1}\widetilde{Q}_{2,h-1} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{h,1}\widetilde{Q}_{h,1}S^{1} & \lambda_{h,2}\widetilde{Q}_{h,2} & \cdots & \lambda_{h,h-1}\widetilde{Q}_{h,h-1} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{1} \\ \widetilde{\mathbf{x}}_{2} \\ \vdots \\ \widetilde{\mathbf{x}}_{h-1} \end{pmatrix} (15)$$

The matrix in (15) — call it M — is a square matrix of dimensions $h\alpha/2 \times h\alpha/2$. A sufficient condition to recover \mathbf{x}_1 is that M is invertible. Notice that the determinant of M, $\det(M)$, is a polynomial in the following variables: $\lambda_{i,j}$, $i \in [h]$, $j \in [h-1]$. Hence, $\det(M)$ is

a nonzero polynomial of degree $h\alpha/2$ in the given variables. From Schwartz–Zippel–DeMillo–Lipton lemma, if the finite field \mathbb{F}_q over which the determinant is defined has cardinality $|\mathbb{F}_q|=q>h\alpha/2$, there exist $\lambda_{i,j}$'s for which the determinant $\det(M)$ above is nonzero.

Part 2: Notice that M above is defined for a particular example scenario. In general, let the number of helper nodes be d, the failed systematic node be $f \in [k]$, the set of helper nodes by $\mathcal{D} \subseteq [n] \setminus \{f\}$, the set of systematic helper nodes be $\mathcal{D}_s \subseteq [k] \setminus \{f\}$, and the set of parity helper nodes be $\mathcal{D}_p \subseteq [k+1:k+r]$. Let the number of parity helper nodes be denoted by h, where h ranges from d-k+1 to r. Let us represent by \mathcal{H}_p the set of parity helper nodes but indexed within [r], where i corresponds to node k+i of the system, that is, $\mathcal{H}_p = \{i \mid k+i \in \mathcal{D}_p\} \subseteq [r]$.

The matrix M in (15), in general, can be seen to be a square matrix of dimensions $h\alpha/(d-k+1) \times h\alpha/(d-k+1)$. In particular, M is a function of f, \mathcal{D}_s , and \mathcal{H}_p , and the determinant polynomial has degree which is a function of $|\mathcal{D}_p| = h$ and d. For each f, \mathcal{D}_s and \mathcal{H}_p , we obtain a sufficiency condition that the corresponding M is invertible. Therefore, the product of the corresponding determinant polynomials is a nonzero polynomial of degree

$$q_{\text{ANY}} = k \left(\sum_{h=d-k+1}^{r} \binom{r}{h} \binom{k-1}{d-h} \frac{h\alpha}{d-k+1} \right)$$
$$= \left(\sum_{h=d-k+1}^{r} h \binom{r}{h} \binom{k-1}{d-h} \right) \frac{k\alpha}{d-k+1};$$

consequently, there exist $\lambda_{i,j}$'s in \mathbb{F} such that any systematic node is repairable with optimal repair bandwidth using any arbitrary set of d helper nodes, as long as the field size $|\mathbb{F}| > q_{\text{ANY}}$.

Part 3: Finally, using Lemma 4, Lemma 5, and Lemma 6, we obtain an [n,k,d] systematic–repair MSR code for any set of d helper nodes, when the field size $q > q_{\text{ANY}} + q_{\text{MDS}}$.

V. Conclusion

In this paper we presented a new construction for systematic–repair MSR codes for all possible values of parameters [n, k, d].

A more generalized construction, where a single [n,k] code simultaneously satisfies the optimal repair for all $d \in \{k+1, \cdots, n-1\}$ will be introduced in a sequel paper. It is to be noted that both these generalizations come at the cost of increasing α . A lower bound on α is proved in [19] when d = n - 1. Whether similar bounds exist for general [n,k,d] or not is left for future work. So is the question of constructing MSR codes that also optimally repair parity nodes.

REFERENCES

[1] D. Borthakur, "The Hadoop Distributed File System: Architecture and Design," in *hadoop.apache.org*, 2007. [Online]. Available: http://hadoop.apache.org/docs/ro.18.0/hdfs_design.pdf

- [2] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google File System," in Proceedings of the 19th ACM Symposium on Operating Systems Principles, 2003, pp. 20-43.
- [3] V. Guruswami and M. Wootters, "Repairing Reed-Solomon Codes," arXiv:1509.04764, 2015.
- [4] A. G. Dimakis, P. B. Godfrey, Y. Wu, M. Wainwright, and K. Ramchandran, "Network Coding for Distributed Storage Systems," in IEEE Transactions on Information Theory, vol. 56, no. 9, September 2010, pp. 4539-4551.
- [5] N. Shah, K. Rashmi, P. V. Kumar, and K. Ramchandran, "Distributed Storage Codes with Repair-by-Transfer and Nonachievability of Interior Points on the Storage-Bandwidth Tradeoff," in IEEE Transactions on Information Theory, vol. 58, no. 3, March 2012, pp. 1837-1852.
- [6] 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," vol. 57, Aug. 2011, pp.
- [7] Y. Wu and A. G. Dimakis, "Reducing Repair Traffic for Erasure Coding-Based Storage via Interference Alignment," in Proceedings of the IEEE International Symposium on Information Theory (ISIT), June-July 2009, pp. 2276–2280.
- [8] D. Cullina, A. G. Dimakis, and T. Ho, "Searching for Minimum Storage Regenerating Codes," in arxiv.org, October 2009. [Online]. Available: http://arxiv.org/abs/0910.2245
- Y. Wu, "A Construction of Systematic MDS Codes with Minimum Repair Bandwidth," in IEEE Transactions on Information Theory, June 2011, pp. 3738-3741.
- [10] V. Cadambe, S. Jafar, H. Maleki, K. Ramchandran, and C. Suh, "Asymptotic Interference Alignment for Optimal Repair of MDS
- codes in Distributed Storage," vol. 59, May 2013, pp. 2974–2987.

 [11] V. R. Cadambe, C. Huang, J. Li, and S. Mehrotra, "Polynomial Length MDS Codes with Optimal Repair in Distributed Storage," in Proceedings of the 45th Asilomar Conference on Signals, Systems and Computers, Nov. 2011, pp. 1850-1854.
- [12] D. S. Papailiopoulos, A. G. Dimakis, and V. R. Cadambe, "Repair Optimal Erasure Codes through Hadamard Designs," in Proceedings of the 49th Annual Allerton Conference on Communication, Con-
- trol, and Computing (Allerton), Sep. 2011, pp. 1382–1389.
 [13] Z. Wang, I. Tamo, and J. Bruck, "On Codes for Optimal Rebuilding Access," in Proceedings of the 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton), September 2011,
- -, "Long MDS Codes for Optimal Repair Bandwidth," in Proceedings of IEEE International Symposium on Information Theory (ISIT), Jul. 2012, pp. 1182-1186.
- [15] C. Suh and K. Ramchandran, "Exact-Repair MDS Codes for Distributed Storage using Interference Alignment," in Proceedings of IEEE International Symposium on Information Theory (ISIT), June 2010, pp. 161-165.
- [16] I. Tamo, Z. Wang, and J. Bruck, "MDS Array Codes with Optimal
- Rebuilding," arXiv:1103.3737v1, 2011.

 [17] I. Tamo and E. En Gad, "[n,k] Minimum-Storage Regenerating Codes for all d's in the Range $\{k, k+1, \ldots, n-1\}$ Simultaneously, 53rd Allerton Conference on Communications, Control, and Computing (invited), 2015.
- [18] A. S. Rawat, O. O. Koyluoglu, and S. Vishwanath, "Progress on High-rate MSR Codes: Enabling Arbitrary Number of Helper Nodes," in arxiv.org, January 2016.
- [19] S. Goparaju, I. Tamo, and R. Calderbank, "An Improved Sub-Packetization Bound for Minimum Storage Regenerating Codes," in IEEE Transactions on Information Theory, vol. 60, no. 5, 2014, pp. 2770-2779.

Appendix A Interference Alignment

Proof of Lemma 1: We prove the result for the failure of systematic node i = 1. The argument generalizes for the failure of other systematic nodes. Let us assume that node x_1 fails, and let each of the remaining d = n - 1 nodes send $\beta = \alpha / r$ symbols to recover \mathbf{x}_1 . In other words, node \mathbf{x}_j (where $j \in [n]$, $j \neq i$) sends $S_i^1 \mathbf{x}_j$ for some $\alpha/r \times \alpha$ matrix S_i^1 . We therefore need to recover \mathbf{x}_1 from the following functions of \mathbf{x}_i , $i \in [k]$:

$$\begin{pmatrix} x_1 & x_2 & x_3 & \cdots & x_k \\ & S_2^1 & & & & \\ & & S_3^1 & & & \\ & & & S_3^1 & & & \\ & & & & \ddots & & \\ S_{k+1}^1 A_{1,1} & S_{k+1}^1 A_{1,2} & \cdots & \cdots & S_{k+1}^1 A_{1,k} \\ S_{k+2}^1 A_{2,1} & S_{k+2}^1 A_{2,2} & \cdots & \cdots & S_{k+2}^1 A_{2,k} \\ \vdots & & \vdots & & \vdots & \vdots & \vdots \\ S_{k+r}^1 A_{r,1} & S_{k+r}^1 A_{r,2} & \cdots & \cdots & S_{k+r}^1 A_{r,k} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_k \end{pmatrix} . (16)$$

Necessity: Suppose the systematic vectors \mathbf{x}_2 through \mathbf{x}_k be the zero vectors. Then, (16) simplifies to:

$$\begin{pmatrix} S_{k+1}^{1}A_{1,1} \\ S_{k+2}^{1}A_{2,1} \\ \vdots \\ S_{k+r}^{1}A_{r,1} \end{pmatrix} \mathbf{x}_{1},$$

where the matrix is an $\alpha \times \alpha$ square matrix. Since \mathbf{x}_1 is recoverable, it is necessary that the matrix be non-singular, thus proving the signal recovery conditions (5) and (7). Note that this also implies that all encoding matrices $A_{i,j}, i \in [r], j \in [k]$, are non-singular. s

Suppose now, without loss of generality, that the interference alignment condition (6) is not satisfied for j = 2. Again, without loss of generality, let $S_2^1 \not\simeq S_{k+1}^1 A_{1,2}$. This implies that

$$\operatorname{rank}\left(\left(\begin{array}{c}S_{2}^{1}\\S_{k+1}^{1}A_{1,2}\end{array}\right)\right) = \frac{\alpha}{r} + \epsilon, \tag{17}$$

for some $\epsilon > 0$. Since \mathbf{x}_1 is recoverable, from (16), we have access to the following information at the replacement node:

From (17), the rank of the matrix in (18) is at least

$$\alpha + \frac{\alpha}{r} + \epsilon + (k-2)\frac{\alpha}{r} > (n-1)\frac{\alpha}{r}$$

the total number of symbols available at the replacement node. In other words, we are able to recover more number of linearly independent symbols that are functions of the systematic data vectors \mathbf{x}_1 through \mathbf{x}_k , than the number of repair symbols available at the replacement node — a contradiction! Thus, conditions (6) and (8) must be true.

Sufficiency: Suppose that we have the required repair matrices S_j^1 that satisfy the signal recovery and interference alignment conditions (5) and (6). Using (6), we can eliminate the contribution of systematic vectors \mathbf{x}_2 through \mathbf{x}_k in the information transmitted by the parity nodes (that is, the last r rows in (16)). For instance, $S_2^1 \subseteq S_{k+1}^1 A_{1,2}$ implies that $S_{k+1}^1 A_{1,2} = BS_2^1$, for some $\alpha/r \times \alpha/r$ matrix B, and therefore the contribution of $S_{k+1}^1 A_{1,2} \mathbf{x}_2$ can be removed from the repair information transmitted by the parity node k+1 using the repair information $S_2^1 \mathbf{x}_2$ (or equivalently, $BS_2^1 \mathbf{x}_2$) transmitted by systematic node 2. Using (5), it is then easy to recover \mathbf{x}_1 .