

Practical Functional Regenerating Codes for Broadcast Repair of Multiple Nodes

Nitish Mital*, Katina Kravevska[†], Cong Ling*, and Deniz Gündüz*

*Department of Electrical Electronics Engineering, Imperial College London

[†]Dep. of Information Security and Communication Technology, NTNU, Norwegian University of Science and Technology
Email: {n.mital,d.gunduz,c.ling}@imperial.ac.uk, katinak@ntnu.no

Abstract—A code construction and repair scheme for optimal functional regeneration of multiple node failures is presented, which is based on stitching together short MDS codes on carefully chosen sets of points lying on a linearized polynomial. The nodes are connected wirelessly, hence all transmissions by helper nodes during a repair round are available to all the nodes being repaired. The scheme is simple and practical because of low subpacketization, low I/O cost and low computational cost. Achievability of the minimum-bandwidth regenerating (MBR) point, as well as an interior point, on the optimal storage-repair bandwidth tradeoff curve is shown. The subspace properties derived in the paper provide insight into the general properties of functional regenerating codes.

I. INTRODUCTION

The content of a file is typically distributed among multiple access points such that accessing any k distinct access points is sufficient to recover the original file. MDS codes provide high storage efficiency while satisfying the above property. When some nodes fail, their cache contents need to be regenerated to be able to continue serving users. An important objective of edge caching in wireless networks is to reduce the backhaul link loads; therefore, we will consider *cache recovery at the edge*; that is, rather than updating the failed cache contents from a central server through backhaul links, the failed cache contents are regenerated with the help of surviving cache nodes. The total amount of data transferred from the surviving nodes to repair the failed nodes is called the *repair bandwidth*. Traditional MDS codes have high storage efficiency, but their repair bandwidth is large [1]. The data of one node is repaired by accessing and transferring data from k nodes, i.e., by recovering the whole content library.

Dimakis et al. showed in [1] that there is a fundamental trade-off between the storage and repair bandwidth by mapping the repair problem in a distributed storage system to a multicasting problem over an information flow graph. The analysis focuses on a single node repair; that is, losing one of the nodes triggers the repair process. Regenerating codes achieve any point on the optimal trade-off curve, while minimum-storage regenerating (MSR) codes and minimum-bandwidth regenerating (MBR) codes operate on the two extremes of this trade-off curve.

This work was supported in part by the European Union's H2020 research and innovation programme under the Marie Skłodowska-Curie Action SCAV-ENGE (grant agreement no. 675891), and by the European Research Council (ERC) Starting Grant BEACON (grant agreement no. 725731).

It was observed in [2] that multiple node repair; that is, the repair process starts only after r nodes fail, is more efficient in terms of the repair bandwidth per node, compared to repairing each node as it fails. In [3] and [4], the authors introduce cooperative regenerating codes, which repair multiple failures cooperatively by allowing each of the r nodes being repaired to collect data from d non-failed nodes, called *helper nodes*, and then to cooperate with the other $r - 1$ nodes being repaired, called *newcomers*. Cooperative repair allows each newcomer to contact any set of helper nodes independently. An explicit construction of regenerating codes that achieve minimum repair bandwidth under cooperative repair is given in [5].

Another model studied in the literature is the centralized repair model [6], [7], in which all the helper nodes transmit the repairing symbols to a centralized node, which then repairs the failed nodes. There being a centralized entity repairing the failed nodes, there is no need for the newcomers to exchange data between themselves like in cooperative repair, thus making the system simpler.

Instead, similarly to [8], we will consider broadcast repair; that is, transmissions from each helper node are received in an error-free manner by all the newcomers. In summary, we will study the broadcast repair of multiple failed cache nodes. The storage-repair bandwidth trade-off for the repair of multiple fully failed nodes is investigated in [9], [10].

The broadcast repair model is theoretically equivalent to the centralized multi-node repair model studied in [6], [7]. Hence, the results that we derive, and codes that we construct for broadcast repair are directly applicable to the centralized repair model. The broadcast repair is different from centralized repair in that, while centralized repair involves a centralized entity which then repairs the failed nodes, thus involving two rounds, broadcast repair involves only one round, and thus, is simpler and faster.

In [7], it is shown that the functional MBR point for repair of multiple nodes is not achievable under exact repair. Similarly, it is shown that under exact repair, the functional repair tradeoff interior points are also not achievable. In [6], [11], it is shown that cooperative repair achieves the minimum bandwidth of centralized repair (and broadcast repair), under exact repair, albeit at a slightly higher storage cost.

In this paper, our contribution is to give an explicit code construction to achieve the optimal storage-repair bandwidth

tradeoff under functional repair, first for the MBR point for all admissible parameters, and then for an interior point, thus potentially providing us with a general framework to construct functional repair codes to achieve any point on the tradeoff curve. The broadcast nature of the system allows all the newcomers to receive the same data, which simplifies the coding scheme. Reference [12] studies the functional repair problem from a projective geometry viewpoint. The subspace/projective geometry view makes it hard to visualize how a set of nodes look like, and how one might approach the construction of a code. There have been attempts to derive the conditions necessary for a functional repair code in [12], [13]. Our code construction strives to address this problem by proposing a simple scheme and simple conditions to guarantee optimality.

Our code construction has the property that for most failure patterns, the helper nodes do not have to perform computations in a repair round; instead they just read and send the data to the newcomers. This property is called repair-by-transfer, which is desirable for a low I/O cost. Functional repair-by-transfer MDS codes were constructed in [14] for some parameters.

Other than these, in our knowledge, there has not been much progress in providing simple, explicit code constructions for functional repair. The network coding literature usually employs random coding, which is not the most practically feasible scheme for distributed storage because of large overheads.

II. SYSTEM MODEL

Consider a wireless caching system where n nodes, each with storage capacity α bits, store a file of size M bits. We index these storage nodes by the set $\mathcal{N} \triangleq \{1, \dots, n\}$. The nodes are fully connected by a wireless broadcast medium and use orthogonal channels for data transmission.

We refer to the nodes that fail as the *failed nodes* and the nodes that do not experience any losses as the *surviving nodes*. We assume that the repair occurs in rounds, where a repair round gets initiated when r nodes experience failures. Thus, a single repair round repairs r failed nodes. There is no loss during a repair round. During a repair round, the failed nodes are repaired with the help of bits transmitted from d surviving nodes, called the *helper nodes*.

A data collector (DC) corresponds to a request to reconstruct the file. Data collectors connect to any subset of k active nodes and retrieve all the stored data in these nodes. This is called the *reconstructability property*. In general, the repair is functional, i.e., the repaired content of the node may not be the same as the original content, but it satisfies the reconstructability property.

Definition 1. The repair bandwidth $\gamma = d\beta$ is defined as the total number of bits the helper nodes broadcast in a repair round.

A. Subspace view

Consider a node storing α linearly independent elements y_1, \dots, y_α from $GF(q^m)$ (possible only if $\alpha \leq m$), then any linear operations performed on these finite field elements, which can be viewed as m -dimensional vectors over $GF(q)$,

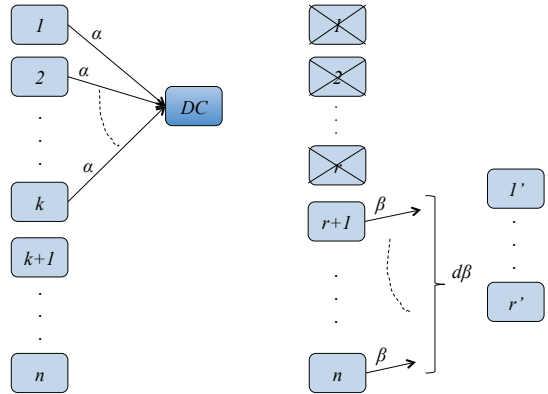


Fig. 1: An illustration of the data collection (reconstructability property) and functional repair of r nodes where d nodes are helpers.

lie in the same subspace. Hence, the node is said to store the subspace of dimension α , denoted by $W_i \equiv \text{span}\{y_i\}, i = 1, \dots, \alpha$. For a set of nodes denoted by \mathcal{A} , the subspace stored by \mathcal{A} is denoted by $W_{\mathcal{A}} = \sum_{i \in \mathcal{A}} W_i$, where the addition operation on subspaces denotes the direct sum of subspaces. By abuse of notation, W_i can also denote the random variables of the stored information in the nodes. The Shannon entropic measures on random variables can be redefined as a measure “ $\dim(\cdot)$ ” (subspace dimension) on intersection, union and set difference of subspaces. The following identities hold [15]:

$$\begin{aligned} H(W_{\mathcal{A}}) &= \dim(W_{\mathcal{A}}) \\ H(W_1, \dots, W_l) &= \dim\left(\sum_{i=1}^l W_i\right) \\ H(W_{\mathcal{A}}|W_{\mathcal{B}}) &= \dim(W_{\mathcal{A}} \setminus W_{\mathcal{B}}) \\ I(W_{\mathcal{A}}; W_{\mathcal{B}}) &= \dim(W_{\mathcal{A}} \cap W_{\mathcal{B}}) \\ I(W_{\mathcal{A}}; W_{\mathcal{B}}; W_{\mathcal{C}}) &= \dim(W_{\mathcal{A}} \cap W_{\mathcal{B}} \cap W_{\mathcal{C}}). \end{aligned}$$

For each set of parameters, a $(n, k, \gamma, d, \alpha, r)$ tuple is feasible, if a code with storage α and repair bandwidth γ exists.

III. MBR POINT CONSTRUCTION FOR MULTIPLE NODE FAILURES

Theorem 1. [16] For any $\alpha \geq \alpha^*(n, k, \gamma, d, r)$, the points $(n, k, \gamma, d, \alpha, r)$ are feasible, and linear network codes suffice to achieve them. If r divides k , the minimum repair bandwidth point is achieved by the pair $(\alpha_{MBR}, \gamma_{MBR}^*) = \frac{2M}{k(2d-k+r)}(d, rd)$.

We provide a construction of a functional repair code for the broadcast setting using linearized polynomials, which were first used by Gabidulin for the construction of rank-metric codes [17].

A. Linearized Polynomials

An important component of our construction is linearized polynomial and their special properties. A linearized polynomial

$$f(x) = \sum_{i=1}^P a_i x^{q^{i-1}}, \quad a_i \in \mathbb{F}_{q^m} \quad (1)$$

can be uniquely identified from evaluations at any P points $x = \theta_i \in \mathbb{F}_{q^m}, i = 1, 2, \dots, P$, that are linearly independent over \mathbb{F}_q . Another relevant property of linearized polynomials is that they satisfy the following condition

$$f(ax + by) = af(x) + bf(y), \quad a, b \in \mathbb{F}_q, x, y \in \mathbb{F}_{q^m}, \quad (2)$$

that is, given a set of points on a linearized polynomial, any linear combination of the points also lies on the polynomial.

B. Code construction

Consider a file of M bits. We split the file into $\frac{k}{2}(2d-k+r)$ packets, denoted by $\{m_1, \dots, m_{\frac{k}{2}(2d-k+r)}\}$. Thus each packet is of size $\frac{2M}{k(2d-k+r)}$ bits. Define the linearized polynomial

$$f(x) = \sum_{i=1}^{\frac{k}{2}(2d-k+r)} m_i x^{q^{i-1}}, \quad m_i \in \mathbb{F}_{q^m} \quad (3)$$

in a field $GF(q^m)$. If a DC receives any $\frac{k}{2}(2d-k+r)$ linearly independent points on the polynomial $f(x)$, it can reconstruct $f(x)$ by interpolation, and thus reconstruct the file.

Pick d linearly independent points on $f(x)$, denoted by $(x_1, y_1), \dots, (x_d, y_d)$, where $y_i = f(x_i)$ for all $i = 1, \dots, d$. Store these points at node 1. Pick another d linearly independent points that are also linearly independent from the points stored at node 1, denoted by $(x_{d+1}, y_{d+1}), \dots, (x_{2d}, y_{2d})$, on $f(x)$, and store them at node 2, and so on, till d nodes are filled with linearly independent points. For the remaining $n-d$ nodes, fill them as if they are being repaired by the d nodes already filled. The scheme for repair is detailed in the Section III-C.

At each node, encode the points with a $(n-1, d)$ systematic MDS code to obtain $n-1$ coded symbols. The symbols at node i are denoted by $\{w_{ij}\}$, where $j \in [n] \setminus \{i\}$.

C. Repair

Suppose the first repair round repairs the nodes $n-r+1$ to n with the helper nodes $1, \dots, d$. Each node $i \in [d]$ transmits the points $w_{ij}, j \in [n-r+1 : n]$. The total number of points broadcasted by d helper nodes is rd . The broadcasted points must be linearly independent, which will become evident that it holds as the repair scheme is explained. Arrange these received points in a $d \times r$ matrix, where row i contains the points transmitted by node i . In the following matrix representation of the received points, we denote the points with the node from which it was transmitted only. It

must be noted that all elements in the following matrix actually represent distinct points.

$$\mathbf{Y} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 2 & 2 & \cdots & 2 \\ \vdots & \ddots & \ddots & \vdots \\ d & d & \cdots & d \end{bmatrix} \quad (4)$$

We then permute the columns of the above matrix in such a way that no row has two points transmitted by the same node. A circular permutation achieving this condition looks like the following

$$\mathbf{Y} = \begin{bmatrix} 1 & 2 & \cdots & r \\ 2 & 3 & \cdots & r+1 \\ \vdots & \ddots & \ddots & \vdots \\ d & 1 & \cdots & r-1 \end{bmatrix}. \quad (5)$$

Now, assume that each node is given a $(2r, r)$ systematic MDS code generator matrix $\mathbf{G} = [\mathbf{I} \ \mathbf{M}]$. \mathbf{G} has rank r , implying that any r columns are linearly independent. Let each newcomer multiply the permuted point-matrix \mathbf{Y} with the local $r \times r$ invertible encoding matrix \mathbf{M} to get $\mathbf{Y}' = \mathbf{Y}\mathbf{M}$. The points in column i of \mathbf{Y}' are stored on the i^{th} newcomer, and encoded with the local $(n-1, d)$ systematic MDS code. The above scheme ensures the following properties to hold:

- L1:** For $i \neq j$, $\dim(W_i \cap W_j) = 0$, or, $I(W_i; W_j) = 0$.
- L2:** For any set of nodes \mathcal{A} , $|\mathcal{A}| \leq r$, the following holds : $\dim(\sum_{i \in \mathcal{A}} W_i) = \sum_{i \in \mathcal{A}} \dim(W_i)$. This is equivalent to the condition - $H(W_{\mathcal{A}}) = \sum_{i \in \mathcal{A}} H(W_i)$. This property is the consequence of permuting the \mathbf{Y} matrix so that packets from the same node are not repeated in the same row. Encoding each row with the full rank matrix \mathbf{M} , extracted from an MDS code, ensures the independence of any r nodes.
- L3:** Given a node A , and a set of nodes denoted by \mathcal{B} , $|\mathcal{B}| \leq d$, partition \mathcal{B} into two disjoint subsets \mathcal{B}_1 and \mathcal{B}_2 . Then the following holds - $\dim(W_A \cap W_{\mathcal{B}_1} \cap W_{\mathcal{B}_2}) = 0$, or, $I(W_A; W_{\mathcal{B}_1}; W_{\mathcal{B}_2}) = 0$. This is because each node transmits a distinct point to repair any node, of which any d of them are linearly independent because of the rank d MDS encoding in each node. Since $d \geq k$, this allows for all admissible parameters.
- L4:** Given a node A , and a set of different r nodes denoted by \mathcal{R} , then $I(W_A; W_{\mathcal{R}}) \leq r$, where equality holds iff 1) the set \mathcal{R} were helper nodes while repairing node A , and a particular row of \mathbf{Y} consisted of the nodes in \mathcal{R} ; or, 2) node A was a helper node while repairing the set \mathcal{R} which failed together.

D. Reconstruction

Suppose a DC accesses the nodes $1, \dots, k$, denoted by \mathcal{A}_{dc} . The points available at the k nodes should be enough to interpolate the polynomial $f(x)$; hence, the necessary and sufficient condition for successful reconstruction is $\dim(W_{\mathcal{A}_{dc}}) \geq$

$\frac{k}{2}(2d-k+r)$. The following lemma will be helpful in showing that the reconstructability property is satisfied.

Lemma 1. Consider a node A , and a set of other $l \leq d$ nodes. Partition the l nodes into sets of r nodes denoted by $\mathcal{R}_1, \dots, \mathcal{R}_{\lfloor l/r \rfloor}$, and the remaining set of nodes denoted by \mathcal{R}' . Then,

$$\dim(W_A \cap \sum_{i=1}^l W_i) = \sum_{i=1}^{\lfloor l/r \rfloor} \dim(W_A \cap W_{\mathcal{R}_i}) \quad (6)$$

Proof.

$$\begin{aligned} \dim(W_A \cap \sum_{i=1}^l W_i) &= I(W_A; W_1, \dots, W_l) \\ &= I(W_A; W_{\mathcal{R}_1}, \dots, W_{\mathcal{R}_{\lfloor l/r \rfloor}}, W_{\mathcal{R}'}) \\ &= I(W_A; W_{\mathcal{R}_1}) + I(W_A; W_{\mathcal{R}_2}, \dots, W_{\mathcal{R}'}, W_{\mathcal{R}_1}) \\ &\stackrel{(a)}{=} I(W_A; W_{\mathcal{R}_1}) + I(W_A; W_{\mathcal{R}_2}, \dots, W_{\mathcal{R}'}) \end{aligned}$$

where (a) holds due to **L3**, which applied to a well known identity from multivariate mutual information, gives $I(W_A; W_{\mathcal{R}_2}, \dots, W_{\mathcal{R}'}, W_{\mathcal{R}_1}) = I(W_A; W_{\mathcal{R}_2}, \dots, W_{\mathcal{R}'}) - I(W_A; W_{\mathcal{R}_2}, \dots, W_{\mathcal{R}'}, W_{\mathcal{R}_1}) = 0$. Using this result inductively, we get

$$\begin{aligned} I(W_A; W_1, \dots, W_l) &= \sum_{i=1}^{\lfloor l/r \rfloor} I(W_A; W_{\mathcal{R}_i}) + I(W_A; W_{\mathcal{R}'}) \\ &\stackrel{(b)}{=} \sum_{i=1}^{\lfloor l/r \rfloor} I(W_A; W_{\mathcal{R}_i}) \\ &= \sum_{i=1}^{\lfloor l/r \rfloor} \dim(W_A \cap W_{\mathcal{R}_i}) \end{aligned} \quad (7)$$

where (b) follows from **L2**. \square

Theorem 2. If a DC accesses k nodes, the dimension of the space obtained from those nodes is $\frac{k}{2}(2d-k+r)$. Thus, the DC can reconstruct the file.

Proof. Without loss of generality, assume that the DC accesses nodes $1, \dots, k$. We have

$$\begin{aligned} \dim\left(\sum_{i=1}^k W_i\right) &= H(W_1, \dots, W_k) \\ &= \sum_{i=1}^k H(W_i | W_{i-1}, \dots, W_1) \\ &= \sum_{i=1}^k [H(W_i) - I(W_i; W_{i-1}, \dots, W_1)] \\ &\stackrel{(c)}{=} \sum_{i=1}^k H(W_i) - \sum_{i=1}^k \sum_{j=1}^{(i-1)/r} I(W_i; W_{\mathcal{R}_j}) \\ &\stackrel{(d)}{\geq} kd - (r(r) + r(2r) + \dots + r(k-r)) \\ &= \frac{k}{2}(2d-k+r) \end{aligned}$$

where (c) follows from Lemma 1, and (d) holds due to **L4**. \square

IV. CODE CONSTRUCTION FOR INTERIOR POINT, $d = n - r$

Theorem 3. [10] The pair $(\alpha, \gamma^*) = \frac{2M}{k(2(n-2r)-(k-r))+2r(k-r)}((n-2r), r(n-r))$ is an interior point on the tradeoff curve, and is achievable with our coding framework.

We consider $d = n - r$ in this section. Each file is divided into $1/2(k(2(n-2r)-(k-r))+2r(k-r))$ subpackets, and the polynomial $f(x)$ is constructed accordingly. The storage and repair scheme is the same as in Section III-C, except that the matrix \mathbf{Y} is of dimensions $(n-2r) \times 2r$, and \mathbf{M} is a $2r \times r$ local encoding matrix, such that there is a $(3r, 2r)$ MDS systematic generator matrix of the form $[\mathbf{I} \ \mathbf{M}]$. Each node stores $n-2r$ linearly independent points, which are then encoded with a local systematic $(n-1, n-2r)$ MDS code. The matrix \mathbf{Y} is constructed by arranging the points received from $n-2r$ helper nodes in the first r columns, like in Equation (4), and the points received from the remaining r helper nodes in the next r columns. The elements of \mathbf{Y} are rearranged so that no two points from a helper node lie in the same row, similar to Equation (5), to obtain the following form,

$$\mathbf{Y} = \begin{bmatrix} \mathbf{S}_{(n-2r) \times r} & \begin{bmatrix} \phi_{r \times r} \\ \vdots \\ \phi_{r \times r} \end{bmatrix} \end{bmatrix}_{(n-2r) \times 2r}$$

where the matrix ϕ consists of the r^2 points from the last r helper nodes. Then, the columns of $\mathbf{Y}' = \mathbf{Y}\mathbf{M}$, which is a $(n-2r) \times r$ matrix, are stored on the r nodes respectively.

The property **L1** is satisfied. Properties **L2, L3** and **L4** are modified as follows:

L2: For any set of nodes \mathcal{A} , $|\mathcal{A}| \leq 2r$, the following holds: $\dim(\sum_{i \in \mathcal{A}} W_i) = \sum_{i \in \mathcal{A}} \dim(W_i)$.

This is equivalent to the condition - $H(W_{\mathcal{A}}) = \sum_{i \in \mathcal{A}} H(W_i)$.

L3: Given a node A , and a set of nodes denoted by \mathcal{B} such that $|\mathcal{B}| \leq n - r$. Partition \mathcal{B} into three disjoint sets $\mathcal{B}_1, \mathcal{B}_2$ and \mathcal{B}_3 , such that $|\mathcal{B}_3| = r$, or in other words, $|\mathcal{B}_1 \cup \mathcal{B}_2| \leq n - 2r$. Then the following holds - $I(W_A; W_{\mathcal{B}_1}; W_{\mathcal{B}_2} | W_{\mathcal{B}_3}) = 0$. This is a consequence of the fact that each node encodes its stored symbols with a rank $n - 2r$ MDS code.

L4: Given a node A , and a set of different $2r$ nodes denoted by \mathcal{B} , the following holds - $I(W_A; W_{\mathcal{B}}) \leq 2r$.

Modified Lemma 1: For any node A , and a set of different nodes denoted by \mathcal{B} , $|\mathcal{B}| \leq n - r$, which is partitioned into sets $\mathcal{R}_1, \dots, \mathcal{R}'$, as in Lemma 1, the following holds due to properties **L2, L3** and **L4**:

$$\begin{aligned} I(W_A; W_{\mathcal{B}}) &= \overline{I(W_A; W_{\mathcal{R}_1})} + I(W_A; W_{\mathcal{B} \setminus \mathcal{R}_1} | W_{\mathcal{R}_1}) \\ &\stackrel{(e)}{=} \sum_{j \geq 2} I(W_A; W_{\mathcal{R}_j} | W_{\mathcal{R}_1}) \end{aligned}$$

where (e) follows similarly to Equation (7) with the slight difference of conditioning. Thus, the dimension of the space obtained by a DC accessing any k nodes is

$$\begin{aligned}
 \dim\left(\sum_{i=1}^k W_i\right) &= H(W_1, \dots, W_k) \\
 &= \sum_{i=1}^k [H(W_i) - I(W_i; W_{i-1}, \dots, W_1)] \\
 &\stackrel{(f)}{=} \sum_{i=1}^k H(W_i) - \sum_{i=1}^k \sum_{j=2}^{(i-1)/r} I(W_i; W_{\mathcal{R}_j} | W_{\mathcal{R}_1}) \\
 &\stackrel{(g)}{\geq} k(n-2r) - (r(r) + \dots + r(k-2r)) \\
 &= \frac{k}{2}(2(n-2r) - (k-r)) + r(k-r)
 \end{aligned}$$

where (f) follows from the modified Lemma 1 above, and (g) holds due to L4.

A. Complexity analysis

Each node must store a local $r \times r$ encoding matrix, and $n-r$ domain points of the polynomial. The subpacketization level of the scheme for the MBR point is $k(2d-k+r)$. Thus functional MBR codes can have a reasonable subpacketization while achieving the optimal repair bandwidth and low I/O cost. The finite field operations are in $GF(q^m)$, $m \geq d^2$. This is because while initially storing points in the nodes, the first d nodes must be filled with d linearly independent points, while the content of the remaining nodes can be generated as if they were being repaired by the first d nodes.

The scheme in this paper is partially repair-by-transfer, when the systematic symbols from the stored points are transmitted. When the non-systematic symbols are transmitted, more than one systematic symbol must be read to form the required linear combination. Another advantage of the scheme is that it stitches together short MDS codes to form a larger code, thus each encoding operation can be done sequentially on a small number of elements.

V. CONCLUSIONS

We presented a practical code construction and repair scheme for functional repair of multiple node failures in a broadcast setting, achieving the MBR point as well as an interior point, as indicated in Fig. 2. We leave for future work to use the construction in this paper for achieving other interior points. The conditions on the stored subspaces provide an insight into the general properties of functional regenerating codes.

REFERENCES

- [1] A. G. Dimakis, P. B. Godfrey, Y. Wu, M. J. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems," *IEEE Trans. on Information Theory*, vol. 56, no. 9, pp. 4539–4551, Sept 2010.
- [2] Y. Hu, Y. Xu, X. Wang, C. Zhan, and P. Li, "Cooperative recovery of distributed storage systems from multiple losses with network coding," *IEEE Journal on Selected Areas in Comms.*, vol. 28, no. 2, pp. 268–276, Feb 2010.

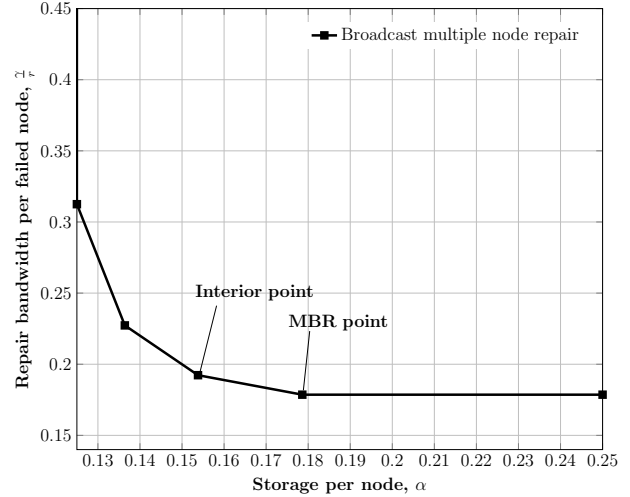


Fig. 2: Repair bandwidth per failed node vs. storage per node for $n = 12$, $k = 8$, $d = 10$, $r = 2$. The points achieved in this paper are illustrated.

- [3] A. M. Kermarrec, N. L. Scouarnec, and G. Straub, "Repairing multiple failures with coordinated and adaptive regenerating codes," in *2011 International Symposium on Networking Coding*, July 2011, pp. 1–6.
- [4] K. W. Shum and Y. Hu, "Cooperative regenerating codes," *IEEE Trans. on Inf. Theory*, vol. 59, no. 11, pp. 7229–7258, Nov 2013.
- [5] A. Wang and Z. Zhang, "Exact cooperative regenerating codes with minimum-repair-bandwidth for distributed storage," in *2013 Proceedings IEEE INFOCOM*, April 2013, pp. 400–404.
- [6] A. S. Rawat, O. O. Koyluoglu, and S. Vishwanath, "Centralized repair of multiple node failures with applications to communication efficient secret sharing," *IEEE Transactions on Information Theory*, vol. 64, no. 12, pp. 7529–7550, Dec 2018.
- [7] M. Zorghi and Z. Wang, "Centralized multi-node repair in distributed storage," in *2016 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, Sept 2016, pp. 617–624.
- [8] M. Gerami, M. Xiao, and M. Skoglund, "Partial repair for wireless caching networks with broadcast channels," *IEEE Wireless Communications Letters*, vol. 4, no. 2, pp. 145–148, April 2015.
- [9] P. Hu, C. W. Sung, and T. H. Chan, "Broadcast repair for wireless distributed storage systems," in *Int. Conf. on Inf., Comms. and Signal Proc.*, Dec 2015, pp. 1–5.
- [10] N. Mital, K. Kravlevska, C. Ling, and D. Gündüz, "Storage-repair bandwidth trade-off for wireless caching with partial failure and broadcast repair," in *2018 IEEE Information Theory Workshop (ITW)*, Nov 2018, pp. 1–5.
- [11] M. Zorghi and Z. Wang, "Centralized multi-node repair regenerating codes," *CoRR*, vol. abs/1706.05431, 2017. [Online]. Available: <http://arxiv.org/abs/1706.05431>
- [12] S. Ng and M. B. Paterson, "Functional repair codes: a view from projective geometry," *CoRR*, vol. abs/1809.08138, 2018. [Online]. Available: <http://arxiv.org/abs/1809.08138>
- [13] H. D. L. Hollmann and W. Poh, "Characterizations and construction methods for linear functional-repair storage codes," in *2013 IEEE International Symposium on Information Theory*, July 2013, pp. 336–340.
- [14] K. W. Shum and Y. Hu, "Functional-repair-by-transfer regenerating codes," in *2012 IEEE International Symposium on Information Theory Proceedings*, July 2012, pp. 1192–1196.
- [15] R. W. Yeung, "A new outlook on shannon's information measures," *IEEE Transactions on Information Theory*, vol. 37, no. 3, pp. 466–474, May 1991.
- [16] A. S. Rawat, O. O. Koyluoglu, and S. Vishwanath, "Centralized repair of multiple node failures," in *IEEE Int. Symp. on Inf. Theory*, July 2016, pp. 1003–1007.
- [17] E. M. Gabidulin, "Theory of codes with maximum rank distance," *Problemy Peredachi Informatsii*, vol. 21, no. 1, pp. 3–16, 1985.