

On the I/O Costs in Repairing Short-Length Reed-Solomon Codes

Weiqi Li*, Hoang Dau†, Zhiying Wang*, Hamid Jafarkhani*, and Emanuele Viterbo‡

*University of California, Irvine, †RMIT University, ‡Monash University

Emails: *{weiqil4, zhiying, hamidj}@uci.edu, †sonhoang.dau@rmit.edu.au, ‡emanuele.viterbo@monash.edu

Abstract—Minimizing the *repair bandwidth*, i.e., the amount of information from the helper nodes needed for recovering the content of one failed node in an erasure-coded distributed storage system, has been the focus of many works in the literature. We investigate another important performance metric, namely the *I/O cost*, which specifies the amount of information that needs to be read by the helper nodes during the repair process of one failed node. We analyze the I/O costs of a few known repair schemes for Reed-Solomon codes of various lengths, in contrast to the previous works in this direction, which only studied the I/O costs in repairing full-length Reed-Solomon codes.

I. INTRODUCTION

Reed-Solomon (RS) codes, although invented back in the sixties [1], still play a crucial role in major distributed storage systems (DSSs), including Google’s Colossus, Quantcast File System, Facebook’s f4, Yahoo Object Store, Baidu’s Atlas, and Backblaze’s Vaults (see [2, Tab. I]). The ubiquity of RS codes stems from their numerous advantages, such as optimal storage overhead, widest range of code parameters, and simple implementations. In the recently released version 3.0.0 of the Hadoop Distributed File System [3], the default erasure coding policy is based on a RS code.

There has been a considerable effort by the research community to optimize the *repair bandwidth* of RS codes [4]–[11]. Here, the repair bandwidth of a repair scheme refers to the amount of data to be transmitted from the helper nodes to the replacement node during the *repair process* that recovers the lost content of one failed node in a distributed storage system deploying the erasure code. Several extensions to the case of multiple erasures were also studied [2], [11]–[14]. The optimal repair bandwidth of RS codes is generally unknown, except for some full-length RS codes [5], [6], [8] and for RS codes with exponential large sub-packetizations [10], [11]. Constructions of RS codes and repair schemes that achieve a trade-off between the repair bandwidth and the sub-packetization were also investigated [15]–[17].

The *I/O cost*, which measures the total amount of data being *read* from the physical disks located at the helper nodes during the repair process of a failed node, is another important performance metric in repairing erasure codes. In the very first work along this line of research [18], the authors showed that known repair schemes of some families of full-length Reed-Solomon codes incur a *trivial I/O cost*, i.e., the entire file needs to be read from the system. Conversely, they also pointed out that for a particular family of RS codes with length $n = 2^\ell$ and two parities over \mathbb{F}_{2^ℓ} , such a high I/O cost is necessary for achieving the optimal repair bandwidth. It was then a natural question to ask whether a trivial I/O cost is always required, or more generally, what is the lowest I/O cost possible. This question was answered for a family of RS codes of length $n = 2^\ell$ and two parities over \mathbb{F}_{2^ℓ} in the subsequent work [19], in which the authors constructed repair schemes that incur an *optimal I/O cost* of $(n-1)\ell - 2^{\ell-1}$ bits.

In this work, we first consider two-parity codes and $q = 2$. We generalize the construction for full-length RS codes [19] to

RS codes of length $n = 2^m$, $m \leq \ell$, and obtain the I/O cost of $(n-1)\ell - 2^{m-1}$ bits. When $m \mid \ell$, we also provide a “lifting” transformation that transforms a repair scheme for a RS code of length q^m over \mathbb{F}_{q^m} with I/O cost I into a repair scheme for a RS code of the same length q^m over the larger field \mathbb{F}_{q^ℓ} with I/O cost $\frac{\ell}{m}I$. Thus we obtain an improved I/O cost of $(n-1)\ell - \frac{\ell}{m}2^{m-1}$ bits. There is still a gap with our newly established lower bound of $(n-1)\ell - (\ell-m+1)2^{m-1}$ bits and it remains an open problem to construct repair schemes achieving optimal I/Os for these RS codes. Second, we consider short-length codes with arbitrary number of parities and field size. We explicitly determine the I/O costs of two existing repair schemes for short-length RS codes established in [10], [16], [17].

II. PRELIMINARIES

Let $[n]$ denote the set $\{1, 2, \dots, n\}$. Let $F = \mathbb{F}_q$ be the finite field of q elements, for some prime power q . Let $E = \mathbb{F}_{q^\ell}$ be an extension field of F , where $\ell \geq 1$, and let $E^* = E \setminus \{0\}$. We refer to the elements of E as *symbols* and the elements of F as *sub-symbols*. The field E may also be viewed as a vector space of dimension ℓ over F , i.e. $E \cong F^\ell$, and hence each symbol in E may be represented as a vector of length ℓ over F . More specifically, suppose $\mathcal{B} = \{\beta_i\}_{i=1}^\ell$ is a basis of E over F , then any element $\alpha \in E$ can be written as $\alpha = \sum_{i=1}^\ell \alpha_i \beta_i$. The unique vector $\phi_{\mathcal{B}}(\alpha) = (\alpha_1, \dots, \alpha_\ell) \in F^\ell$ is called the *vector representation* of α w.r.t. the basis \mathcal{B} .

The (field) trace of a symbol $\alpha \in E$ over F is defined to be $\text{Tr}_{E/F}(\alpha) = \sum_{i=0}^{\ell-1} \alpha^{q^i}$ (the subscript E/F is often omitted). The support of a vector $\mathbf{u} = (u_1, \dots, u_\ell)$, denoted $\text{supp}(\mathbf{u})$, is the set $\{j: u_j \neq 0\}$. The (Hamming) weight of \mathbf{u} , denoted $\text{wt}(\mathbf{u})$, is $|\text{supp}(\mathbf{u})|$. The support of a set of vectors U is $\text{supp}(U) \triangleq \bigcup_{\mathbf{u} \in U} \text{supp}(\mathbf{u})$. A *linear* $[n, k]$ code \mathcal{C} over E is an E -subspace of E^n of dimension k . Each element of a code is referred to as a *codeword*. The *dual* of a code \mathcal{C} , denoted \mathcal{C}^\perp , is the orthogonal complement of \mathcal{C} in E^n and has dimension $r = n - k$. Elements of \mathcal{C}^\perp are called *dual codewords*.

Definition 1. Let $E[x]$ denote the ring of polynomials over a finite field E . The Reed-Solomon code $\text{RS}_E(\mathcal{A}, k) \subseteq E^n$ of dimension k with evaluation points $\mathcal{A} = \{\alpha_j\}_{j=1}^n \subseteq E$ is defined as

$$\text{RS}_E(\mathcal{A}, k) = \left\{ (f(\alpha_1), \dots, f(\alpha_n)) : f \in E[x], \deg(f) < k \right\}.$$

The RS code is *full length* if $n = |E|$. When \mathcal{A} forms an F -subspace of E , the dual of the RS code $\text{RS}_E(\mathcal{A}, k)$ is another RS code $\text{RS}_E(\mathcal{A}, n - k)$ (see Section IV-A).

Trace repair framework. First, note that each symbol in E can be recovered from ℓ independent traces. More precisely, given a basis $\{\beta_i\}_{i=1}^\ell$ of E over F , any $\alpha \in E$ can be uniquely determined given the values of $\text{Tr}(\alpha\beta_i)$ for $i \in [\ell]$, i.e. $\alpha = \sum_{i=1}^\ell \text{Tr}(\alpha\beta_i)\beta'_i$, where $\{\beta'_i\}_{i=1}^\ell$ is the *dual (trace-orthogonal) basis* of $\{\beta_i\}_{i=1}^\ell$ (see, e.g., [20, Ch. 2, Def. 2.30]).

Let \mathcal{C} be an $[n, k]$ linear code over E and \mathcal{C}^\perp be its dual. If $\mathbf{c} = (c_1, \dots, c_n) \in \mathcal{C}$ and $\mathbf{g} = (g_1, \dots, g_n) \in \mathcal{C}^\perp$ then $\mathbf{c} \cdot \mathbf{g}$

$= \sum_{j=1}^n c_j \mathbf{g}_j = 0$. Note that when \mathcal{A} is an F -subspace of E and $\mathcal{C} = \text{RS}_E(\mathcal{A}, k)$, \mathbf{g} is a codeword of $\text{RS}_E(\mathcal{A}, n - k)$ and hence corresponds to a polynomial $g(x)$ of degree less than $n - k$. Suppose c_{j^*} is erased and needs to be recovered. In the trace repair framework, depending on j^* , we choose a set of ℓ dual codewords $\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(\ell)}$ such that $\dim_F(\{\mathbf{g}_{j^*}^{(i)}\}_{i=1}^\ell) = \ell$. Since trace is linear, we obtain the following ℓ equations

$$\text{Tr}(\mathbf{g}_{j^*}^{(i)} c_{j^*}) = - \sum_{j \neq j^*} \text{Tr}(\mathbf{g}_j^{(i)} c_j), \quad i \in [\ell]. \quad (1)$$

In order to recover c_{j^*} , one needs to retrieve sufficient information from $\{c_j\}_{j \neq j^*}$ to compute the right-hand sides of (1). We define, for every $j \in [n]$,

$$\mathcal{S}_{j \rightarrow j^*} \triangleq \text{span}_F(\{\mathbf{g}_j^{(1)}, \dots, \mathbf{g}_j^{(\ell)}\}). \quad (2)$$

Then for each $j \neq j^*$, to determine $\text{Tr}(\mathbf{g}_j^{(i)} c_j)$ for all $i \in [\ell]$, it suffices to retrieve $\dim_F(\mathcal{S}_{j \rightarrow j^*})$ sub-symbols (in F) only. Indeed, suppose $\{\mathbf{g}_j^{(i_t)}\}_{t=1}^s$ is an F -basis of $\mathcal{S}_{j \rightarrow j^*}$, then by retrieving just s traces $\text{Tr}(\mathbf{g}_j^{(i_1)} c_j), \dots, \text{Tr}(\mathbf{g}_j^{(i_s)} c_j)$ of c_j , all other traces $\text{Tr}(\mathbf{g}_j^{(i)} c_j)$ can be computed as F -linear combinations of those s traces without any knowledge of c . Finally, since $\{\mathbf{g}_{j^*}^{(i)}\}_{i=1}^\ell$ is F -linearly independent, c_{j^*} can be recovered from its ℓ corresponding traces on the left-hand side of (1). We refer to such a scheme as a *repair scheme based on $\{\mathbf{g}^{(i)}\}_{i=1}^\ell$* . It was known that this type of repair schemes includes every possible linear repair scheme for RS codes [5].

Lemma 1 (Guruswami-Wootters [5]). *Suppose $E = \mathbb{F}_{q^\ell}$, $F = \mathbb{F}_q$, \mathcal{C} is an $[n, k]$ linear code over E and \mathcal{C}^\perp is its dual. The repair scheme for c_{j^*} based on ℓ dual codewords $\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(\ell)}$, where $\dim_F(\mathcal{S}_{j^* \rightarrow j^*}) = \ell$, incurs a repair bandwidth of $\sum_{j \neq j^*} \dim_F(\mathcal{S}_{j \rightarrow j^*})$ sub-symbols in F , where $\mathcal{S}_{j \rightarrow j^*}$ is defined as in (2).*

I/O Cost of a Repair Scheme Let $\mathcal{B} = \{\beta_i\}_{i=1}^\ell$ be an F -basis of E . Then each element $\alpha = \sum_{i=1}^\ell \alpha_i \beta_i \in E$ can be represented by a vector $\phi(\alpha) = (\alpha_1, \dots, \alpha_\ell) \in F^\ell$ as defined earlier. We assume throughout this work that every node uses a fixed (but probably different) basis to represent and store the finite field elements. Another underlying assumption is that each sub-symbol α_i of α can be read from the storage disk *separately* without accessing other sub-symbols. We first define the I/O cost of a function and then proceed to describe the I/O cost of a repair scheme.

Definition 2 ([18]). The (read) I/O cost of a function $f(\cdot)$ w.r.t. a basis \mathcal{B} is the minimum number of sub-symbols of $\alpha \in E$ needed for the computation of $f(\alpha)$. The I/O cost of a set of functions \mathcal{F} is the minimum number of sub-symbols of α needed for the computation of $\{f(\alpha) : f \in \mathcal{F}\}$.

Lemma 2 ([18]). *The following statements hold.*

- (a) *The I/O cost of the trace functional $\text{Tr}^\gamma(\cdot)$, defined by $\text{Tr}^\gamma(\alpha) \triangleq \text{Tr}(\gamma\alpha)$, w.r.t. \mathcal{B} is $\text{wt}(\mathbf{w}^{\gamma, \mathcal{B}})$, where $\mathbf{w}^{\gamma, \mathcal{B}} \triangleq (\text{Tr}(\gamma\beta_1), \dots, \text{Tr}(\gamma\beta_\ell)) \in F^\ell$. (3)*
- (b) *The I/O cost of the set of trace functionals $\{\text{Tr}^\gamma(\cdot) : \gamma \in \Gamma\}$ w.r.t. \mathcal{B} is $|\cup_{\gamma \in \Gamma} \text{supp}(\mathbf{w}^{\gamma, \mathcal{B}})|$.*

The I/O cost of the repair scheme for c_{j^*} based on a set of dual codewords $\{\mathbf{g}^{(i)}\}_{i=1}^\ell$ is the minimum number of sub-symbols of c_j 's, $j \neq j^*$, needed in the computation of the right-hand side of (1). Note that each node may use a fixed but different basis \mathcal{B}_j , $j \in [n]$, to represent the finite field elements. The formal definition is given below.

Definition 3. The I/O cost of the repair scheme for c_{j^*} based on a set of dual codewords $\{\mathbf{g}^{(i)}\}_{i=1}^\ell$ w.r.t. a set of bases

$\{\mathcal{B}_j\}_{j \neq j^*}$ is the sum of the I/O costs of the sets of trace functionals $\mathcal{F}_j = \{\text{Tr}^{\mathbf{g}_j^{(i)}}(\cdot)\}_{i=1}^\ell$ w.r.t. \mathcal{B}_j , $j \in [n] \setminus \{j^*\}$.

Lemma 3 follows directly from Lemma 2 and Definition 3.

Lemma 3. *The I/O cost of the repair scheme for c_{j^*} based on a set of dual codewords $\{\mathbf{g}^{(i)}\}_{i=1}^\ell$ w.r.t. a set of bases $\{\mathcal{B}_j\}_{j \neq j^*}$ is $\sum_{j \in [n] \setminus \{j^*\}} \text{nz}(\mathbf{W}_j)$, where $\text{nz}(\mathbf{W}_j)$ specifies the number of nonzero columns in the $\ell \times \ell$ I/O matrix \mathbf{W}_j defined as*

$$\mathbf{W}_j \triangleq \begin{pmatrix} \mathbf{w}^{\mathbf{g}_j^{(1)}, \mathcal{B}_j} \\ \mathbf{w}^{\mathbf{g}_j^{(2)}, \mathcal{B}_j} \\ \vdots \\ \mathbf{w}^{\mathbf{g}_j^{(\ell)}, \mathcal{B}_j} \end{pmatrix} = \begin{pmatrix} \text{Tr}(\mathbf{g}_j^{(1)} \beta_{1,j}) \cdots \text{Tr}(\mathbf{g}_j^{(1)} \beta_{\ell,j}) \\ \text{Tr}(\mathbf{g}_j^{(2)} \beta_{1,j}) \cdots \text{Tr}(\mathbf{g}_j^{(2)} \beta_{\ell,j}) \\ \vdots \\ \text{Tr}(\mathbf{g}_j^{(\ell)} \beta_{1,j}) \cdots \text{Tr}(\mathbf{g}_j^{(\ell)} \beta_{\ell,j}) \end{pmatrix},$$

where $\mathcal{B}_j = \{\beta_{i,j}\}_{i=1}^\ell$. In this repair scheme the node storing c_j must read the i -th sub-symbol of c_j if and only if the i -th column of \mathbf{W}_j is nonzero. Thus, $\text{nz}(\mathbf{W}_j)$ specifies the I/O cost incurred at the node storing c_j .

Lemma 3 can be understood as follows. If $c_j = f(\alpha_j)$ is read in full, ℓ sub-symbols will need to be accessed. However, it is possible that fewer than ℓ sub-symbols of c_j need to be read. Each zero column of \mathbf{W}_j indicates a sub-symbol in the vector representation of c_j that does not need to be read.

III. ON THE I/O COST FOR $[2^m, 2^m - 2]_{2^\ell}$ RS CODES

In this section we consider $[2^m, 2^m - 2]_{2^\ell}$ RS codes and assume that the nodes use a common F -basis of E . We first establish a lower bound on the I/O cost, thus generalizing [19, Thm. 1]. Then we present one construction generalizing [19, Cnstr. 1], and one based on a lifting transformation.

Theorem 1. *For a Reed-Solomon code $\text{RS}_{\mathbb{F}_{2^\ell}}(\mathcal{S}_m, 2^m - 2)$ where \mathcal{S}_m is an m -dimensional \mathbb{F}_2 -subspace of \mathbb{F}_{2^ℓ} , the I/O cost (in bits) of an arbitrary linear repair scheme \mathcal{R} always satisfies the following inequality.*

$$\text{ic}(\mathcal{R}) \geq (n - 1)\ell - (\ell - m + 1)2^{m-1}. \quad (4)$$

We first prove this theorem and then provide two repair schemes. Although not achieving the lower bound, they may provide a hint for constructing an I/O-optimal one.

As $r = 2$ and \mathcal{S}_m is an m -dimensional subspace of \mathbb{F}_{2^ℓ} , a dual codeword of $\mathcal{C}_m = \text{RS}_{\mathbb{F}_{2^\ell}}(\mathcal{S}_m, 2^m - 2)$ can be obtained by evaluating a polynomial of degree at most one at all the elements of \mathcal{S}_m . A linear repair scheme, therefore, is based on a set of ℓ polynomials $g_i(x) = \mathbf{a}_i x + \mathbf{b}_i$, $\mathbf{a}_i, \mathbf{b}_i \in \mathbb{F}_{2^\ell}$, $i \in [\ell]$. Set $A = \{\mathbf{a}_i\}_{i=1}^\ell$ and $B = \{\mathbf{b}_i\}_{i=1}^\ell$. We say the repair scheme is *defined by A and B* . Moreover, by generalizing [18, Lem. 8] to the case when the evaluation points form a subspace, it suffices to consider repairing $c_1 = f(0)$, the first codeword component corresponding to the evaluation point $\alpha_1 = 0$. As a repair scheme for c_1 , it is required that $\text{rank}_{\mathbb{F}_2}(\{g_i(0)\}_{i=1}^\ell) = \ell$. In other words, B must be an \mathbb{F}_2 -basis of \mathbb{F}_{2^ℓ} . We henceforth set $r_A \triangleq \text{rank}_{\mathbb{F}_2}(A)$ and $A\gamma + B \triangleq \{\mathbf{a}_i \gamma + \mathbf{b}_i\}_{i=1}^\ell$.

It will later become clear in the proof of Theorem 1 that the set of “good” helpers w.r.t. a fixed basis element $\beta = \beta_i \in \mathcal{B}$, defined in Lemma 4, consists of the helpers where one bit of I/O can be saved in the repair scheme defined by A and B .

Lemma 4. *Let \mathcal{S}_m be an m -dimensional \mathbb{F}_2 -subspace of \mathbb{F}_{2^ℓ} . Suppose $A = \{\mathbf{a}_i\}_{i=1}^\ell \subset \mathbb{F}_{2^\ell}$, $B = \{\mathbf{b}_i\}_{i=1}^\ell$ is an \mathbb{F}_2 -basis of \mathbb{F}_{2^ℓ} , and $\beta \in \mathbb{F}_{2^\ell}^*$. We define $G_{A,B,\beta}^m \triangleq \mathcal{S}_m \cap G_{A,B,\beta}$ where*

$$G_{A,B,\beta} \triangleq \{\gamma \in \mathbb{F}_{2^\ell} : \text{Tr}((\mathbf{a}_i \gamma + \mathbf{b}_i)\beta) = 0, \forall i \in [\ell]\}. \quad (5)$$

Then $G_{A,B,\beta}^m$, which is called the set of “good” helpers w.r.t. β , has size at most $2^{\min\{m-1, \ell-r_A\}}$.

Next, we propose another construction based on Proposition 1, which transforms a repair scheme for a full-length code over \mathbb{F}_{q^m} into a repair scheme for the code of the same length $n = q^m$ but over a larger field $\mathbb{F}_{q^\ell} \supset \mathbb{F}_{q^m}$. The proposition applies to arbitrary number of parities and arbitrary q .

Proposition 1 (Lifting transformation). *Let \mathcal{R}_m be a repair scheme of $\mathcal{C}_m = \text{RS}_{\mathbb{F}_{q^m}}(\mathcal{A}, k)$ where $\mathcal{A} = \mathbb{F}_{q^m}$. If $m \mid \ell$ then there exists a repair scheme \mathcal{R}_ℓ of $\mathcal{C}_\ell = \text{RS}_{\mathbb{F}_{q^\ell}}(\mathcal{A}, k)$ so that $\text{ic}(\mathcal{R}_\ell) = \frac{\ell}{m} \text{ic}(\mathcal{R}_m)$, where $\text{ic}(\cdot)$ denotes the I/O cost of a repair scheme w.r.t. an appropriately selected set of bases.*

Proof. Suppose \mathcal{R}_m repairs c_{j^*} of the codeword $c \in \mathcal{C}_m$ and is based on the set of dual codewords $\{g^{(i)}\}_{i=1}^m$ obtained by evaluating the polynomials $\{g^{(i)}(x)\}_{i=1}^m$ of degree less than k at the points of \mathcal{A} . Moreover, suppose $\mathcal{B}^m = \{\beta_i\}_{i=1}^m$ is the common \mathbb{F}_q -basis of \mathbb{F}_{q^m} used by every node in \mathcal{R}_m .

Let $\{\tau_t\}_{t=1}^{\ell/m}$ be an \mathbb{F}_{q^m} -basis of \mathbb{F}_{q^ℓ} and $\{\eta_t\}_{t=1}^{\ell/m}$ its dual (trace-orthogonal) basis. Consider the scheme \mathcal{R}_ℓ that repairs \bar{c}_{j^*} of the codeword $\bar{c} \in \mathcal{C}_\ell$ based on the set of dual codewords $\{\tau_t g^{(i)}\}_{t \in [\ell/m], i \in [m]}$. Note that $\tau_t g^{(i)}$ is obtained by evaluating the polynomial $\tau_t g^{(i)}(x)$ of degree less than k on \mathcal{A} and is, therefore, a dual codeword of $\text{RS}_{\mathbb{F}_{q^\ell}}(\mathcal{A}, k)$. Let $\mathcal{B}^\ell = \{\eta_t \beta_i\}_{t \in [\ell/m], i \in [m]}$ be the common basis used by every node in \mathcal{R}_ℓ . Note that since $\text{rank}_q(\{g_j^{(i)}\}_{i=1}^m) = m$, we have $\text{rank}_q(\{\tau_t g_j^{(i)}\}_{t \in [\ell/m], i \in [m]}) = \ell$, and hence, \mathcal{R}_ℓ is indeed a repair scheme for \bar{c}_{j^*} .

We subsequently show that to repair \bar{c}_{j^*} , the I/O cost incurred at the node storing \bar{c}_j , $j \neq j^*$, w.r.t. \mathcal{B}^ℓ , is ℓ/m times more than the I/O cost incurred at the node storing c_j when repairing c_{j^*} w.r.t. \mathcal{B}^m . To this end, let \mathbf{W}_j^m and \mathbf{W}_j^ℓ denote the I/O matrices corresponding to c_j and \bar{c}_j , respectively. The rows and columns of \mathbf{W}_j^m correspond to the sets $\{g^{(i)}\}_{i=1}^m$ and $\{\beta_i\}_{i=1}^m$, respectively, and its typical (i, i') -entry is $\text{Tr}_{\mathbb{F}_{q^m}/\mathbb{F}_q}(g_j^{(i)} \beta_{i'})$. The rows and columns of \mathbf{W}_j^ℓ correspond to the sets $\{\tau_t g^{(i)}\}_{t \in [\ell/m], i \in [m]}$ and $\{\eta_t \beta_i\}_{t \in [\ell/m], i \in [m]}$, respectively, and its typical $((t, i), (t', i'))$ -entry is

$$\begin{aligned} \text{Tr}_{\mathbb{F}_{q^\ell}/\mathbb{F}_q}(\tau_t g_j^{(i)} \eta_{t'} \beta_{i'}) &= \text{Tr}_{\mathbb{F}_{q^m}/\mathbb{F}_q}(\text{Tr}_{\mathbb{F}_{q^\ell}/\mathbb{F}_{q^m}}(\tau_t g_j^{(i)} \eta_{t'} \beta_{i'})) \\ &= \text{Tr}_{\mathbb{F}_{q^m}/\mathbb{F}_q}(g_j^{(i)} \beta_{i'} \text{Tr}_{\mathbb{F}_{q^\ell}/\mathbb{F}_{q^m}}(\tau_t \eta_{t'})) \\ &= \begin{cases} \text{Tr}_{\mathbb{F}_{q^m}/\mathbb{F}_q}(g_j^{(i)} \beta_{i'}), & \text{if } t = t', \\ 0, & \text{if } t \neq t', \end{cases} \end{aligned}$$

because $g_j^{(i)}, \beta_{i'} \in \mathbb{F}_{q^m}$ and $\{\tau_t\}_{t=1}^{\ell/m}$ and $\{\eta_t\}_{t=1}^{\ell/m}$ are dual bases of \mathbb{F}_{q^ℓ} over \mathbb{F}_{q^m} . Therefore, the $\ell \times \frac{\ell}{m}$ submatrix of \mathbf{W}_j^ℓ consisting of columns indexed by the set $\{(t', i')\}_{t'=1}^{\ell/m}$ can be obtained by taking the Kronecker product of the identity matrix of order ℓ/m and the i' -th column of \mathbf{W}_j^m . Therefore, in Lemma 3's notation, $\text{nz}(\mathbf{W}_j^\ell) = \frac{\ell}{m} \text{nz}(\mathbf{W}_j^m)$, which implies that $\text{ic}(\mathcal{R}_\ell) = \frac{\ell}{m} \text{ic}(\mathcal{R}_m)$. ■

When $m \mid \ell$, the lifting transformation provided in Proposition 1 yields a lower I/O cost than the one in Construction I. However, a specific basis must be used.

Corollary 1. *If $m \mid \ell$ then there exists a repair scheme of the code $\text{RS}_{\mathbb{F}_{2^\ell}}(\mathcal{A}, 2^m - 2)$ where $\mathcal{A} = \mathbb{F}_{2^m} \subset \mathbb{F}_{2^\ell}$ with an I/O cost of $(n - 1)\ell - \frac{\ell}{m} 2^{m-1}$ bits.*

Proof. As shown in [19], there exists a repair scheme \mathcal{R}_m of $\mathcal{C}_m = \text{RS}_{\mathbb{F}_{2^m}}(\mathcal{A}, 2^m - 2)$ where $\mathcal{A} = \mathbb{F}_{2^m}$ with an I/O cost of $(2^m - 1)m - 2^{m-1}$. By Proposition 1, there exists a repair scheme \mathcal{R}_ℓ of $\mathcal{C}_\ell = \text{RS}_{\mathbb{F}_{2^\ell}}(\mathcal{A}, 2^m - 2)$ with the I/O cost

$$\text{ic}(\mathcal{R}_\ell) = \frac{\ell}{m} ((2^m - 1)m - 2^{m-1}) = (n - 1)\ell - \frac{\ell}{m} 2^{m-1}. \quad \blacksquare$$

IV. ON THE I/O COSTS OF OTHER FAMILIES OF RS CODES

In this section we consider arbitrary number of parities. We show that with a careful selection of the basis at each node, known repair schemes for two families of RS codes can achieve low I/O costs. For the first family, the I/O cost equals the repair bandwidth. The second family achieves the optimal repair bandwidth, and we show that the I/O cost is at most twice the bandwidth. We henceforth assume that different nodes may use different bases but known to all other nodes. We use a different but equivalent way to compute the I/O cost.

A. Another Formulation of the I/O Cost

The dual of $\text{RS}_E(\mathcal{A}, k)$ where \mathcal{A} is not necessarily a subspace is a Generalized Reed-Solomon (GRS) code, given by $\{(\mu_\alpha g(\alpha))_{\alpha \in \mathcal{A}} : g \in E[x], \deg(g) < n - k\}$, where μ_α is the (nonzero) column multiplier corresponding to the evaluation point $\alpha \in \mathcal{A}$. Moreover, $\mu_\alpha = 1 / \prod_{\alpha' \in \mathcal{A} \setminus \{\alpha\}} (\alpha - \alpha')$ [21, p. 167], which can be different in general. When \mathcal{A} is a subspace of E , $\mu_\alpha = 1 / \prod_{\alpha' \in \mathcal{A}^*} \alpha'$ and are the same for every $\alpha \in \mathcal{A}$, which is the case in the previous sections.

If $E = \mathbb{F}_{q^\ell}$ and $F = \mathbb{F}_q$ then each repair scheme for $c_{\alpha^*} = f(\alpha^*)$ is based on a set of ℓ polynomials $g_i(x) \in E[x]$, $i \in [\ell]$, $\deg(g_i) < r = n - k$, and the corresponding repair equations are

$$\text{Tr}_{E/F}(\mu_{\alpha^*} g_i(\alpha^*) f(\alpha^*)) = - \sum_{\alpha \in \mathcal{A} \setminus \{\alpha^*\}} \text{Tr}_{E/F}(\mu_\alpha g_i(\alpha) f(\alpha)).$$

Let $\mathcal{B}_\alpha = \{\beta_{\alpha, i}\}_{i=1}^\ell$ be the F -basis of E used by the node storing $f(\alpha)$. Note that this basis should be fixed for α and independent of the failed node. By Lemma 2 (c), the computation of $\text{Tr}_{E/F}(\mu_\alpha g_i(\alpha) f(\alpha))$ incurs an I/O cost equal to the Hamming weight of the vector $w^{\mu_\alpha g_i(\alpha), \mathcal{B}_\alpha}$, which is

$$(\text{Tr}_{E/F}(g_i(\alpha) \mu_\alpha \beta_{\alpha, 1}), \dots, \text{Tr}_{E/F}(g_i(\alpha) \mu_\alpha \beta_{\alpha, \ell})).$$

Let $\mathcal{B}'_\alpha = \{\beta'_{\alpha, i}\}_{i=1}^\ell$ be the dual (trace-orthogonal) basis of $\mu_\alpha \mathcal{B}_\alpha = \{\mu_\alpha \beta_{\alpha, i}\}_{i=1}^\ell$. Then, the I/O cost of each of the trace functionals on the RHS of the repair equations above is the Hamming weight of $\phi_{\mathcal{B}'_\alpha}(g_i(\alpha))$, the vector representation of $g_i(\alpha)$ w.r.t. the basis \mathcal{B}'_α . Therefore, the I/O cost incurred at the node storing $f(\alpha)$ is equal to the number of elements in \mathcal{B}'_α required to linearly generate $g_1(\alpha), \dots, g_\ell(\alpha)$ over \mathbb{F}_q . We henceforth use this formulation of the I/O cost.

B. Two-Coset Codes

In this subsection, we slightly modify the construction in [17, Thm. 2] to obtain repair schemes with the I/O costs as low as the repair bandwidths established in their work.

Suppose $m \mid \ell$ and β and γ are the primitive elements of \mathbb{F}_{q^ℓ} and \mathbb{F}_{q^m} , respectively. Consider an even $n \leq 2(q^m - 1)$ and $k > 0$ so that $\ell/m \leq r = n - k$. For $0 < m_1 < m_2$ satisfying $m_2 - m_1 = q^{sm}$ (the reason to have this will be clear later), we consider the so-called two-coset code $\text{RS}_{\mathbb{F}_{q^\ell}}(\mathcal{A}, k)$ where \mathcal{A} consists of $n/2$ points from the coset $\beta^{m_1} \mathbb{F}_{q^m}^*$ and $n/2$ points from the coset $\beta^{m_2} \mathbb{F}_{q^m}^*$ of $\mathbb{F}_{q^m}^*$ in $\mathbb{F}_{q^\ell}^*$. As β is a primitive element of \mathbb{F}_{q^ℓ} , these two cosets are disjoint.

Suppose that the node storing $c_\alpha = f(\alpha)$ uses an \mathbb{F}_q -basis \mathcal{B}_α so that $\mathcal{B}'_\alpha = \{\gamma^{t_1} \beta^{t_2} : 0 \leq t_1 \leq m - 1, 0 \leq t_2 \leq \ell/m - 1\}$, where \mathcal{B}_α and \mathcal{B}'_α are defined in Section IV-A. Note that even if \mathcal{B}'_α does not depend on α , the column multiplier μ_α does, and so does the basis \mathcal{B}_α used by the node. We consider the repair scheme \mathcal{R}_{α^*} for c_{α^*} based on the set of polynomials

$$g_{\alpha^*, (t_1, t_2)}(x) = \begin{cases} \gamma^{t_1} \left(\frac{x}{\beta^{m_2}}\right)^{t_2}, & \text{if } \alpha^* \in \beta^{m_1} \mathbb{F}_{q^m}^*, \\ \gamma^{t_1} \left(\frac{x}{\beta^{m_1}}\right)^{t_2}, & \text{if } \alpha^* \in \beta^{m_2} \mathbb{F}_{q^m}^*, \end{cases}$$

for $0 \leq t_1 \leq m - 1$ and $0 \leq t_2 \leq \ell/m - 1$.

We now analyze the I/O cost of \mathcal{R}_{α^*} . Without loss of generality, we assume that $\alpha^* \in \beta^{m_2} \mathbb{F}_{q^m}^*$. First, for $\alpha \in \beta^{m_1} \mathbb{F}_{q^m}^*$, i.e., $\alpha = \beta^{m_1} \gamma^{t_0}$ for some $0 \leq t_0 \leq q^m - 2$, we have $g_{\alpha^*, (t_1, t_2)}(\alpha) = \gamma^{t_1 + t_0 t_2} \in \mathbb{F}_{q^m}^*$, for every t_1, t_2 . Therefore, $g_{\alpha^*, (t_1, t_2)}(\alpha)$ can always be generated as linear combinations of the m elements $\{\gamma^0, \gamma^1, \dots, \gamma^{m-1}\} \subset \mathcal{B}'_{\alpha}$. Thus, the I/O cost incurred at the node storing c_{α} if $\alpha \in \beta^{m_1} \mathbb{F}_{q^m}^*$ is at most m . Next, suppose that $\alpha = \beta^{m_2} \gamma^{t_0} \in \beta^{m_2} \mathbb{F}_{q^m}^*$. Then, $g_{\alpha^*, (t_1, t_2)}(\alpha) = \beta^{(m_2 - m_1)t_2} \gamma^{t_1 + t_0 t_2}$. As $\{\beta^{t_2}\}_{t_2=0}^{\ell/m-1}$ is an \mathbb{F}_{q^m} -basis of $\mathbb{F}_{q^{\ell}}$, $\{\beta^{(m_2 - m_1)t_2}\}_{t_2=0}^{\ell/m-1} = \{\beta^{q^{sm} t_2}\}_{t_2=0}^{\ell/m-1}$ is also an \mathbb{F}_{q^m} -basis of $\mathbb{F}_{q^{\ell}}$ (cf. [17, Lem. 2]). As $\gamma^{t_0 t_2} \in \mathbb{F}_{q^m}$, this implies that $\{\beta^{(m_2 - m_1)t_2} \gamma^{t_0 t_2}\}_{t_2=0}^{\ell/m-1}$ is also an \mathbb{F}_{q^m} -basis of $\mathbb{F}_{q^{\ell}}$. Moreover, since $\{\gamma^{t_1}\}_{t_1=0}^{m-1}$ is an \mathbb{F}_q -basis of \mathbb{F}_{q^m} , the set $\{g_{\alpha^*, (t_1, t_2)}(\alpha) : 0 \leq t_1 \leq m-1, 0 \leq t_2 \leq \ell/m-1\} = \{\beta^{(m_2 - m_1)t_2} \gamma^{t_0 t_2} \gamma^{t_1} : 0 \leq t_1 \leq m-1, 0 \leq t_2 \leq \ell/m-1\}$ forms an \mathbb{F}_q -basis of $\mathbb{F}_{q^{\ell}}$. Hence, the I/O cost incurred at the node storing c_{α} if $\alpha \in \beta^{m_2} \mathbb{F}_{q^m}^*$ is ℓ . Summing up the two cases,

$$\text{ic}(\mathcal{R}_{\alpha^*}) \leq (n/2)m + (n/2 - 1)\ell,$$

which is the same as the repair bandwidth in [16, Thm. 2].

C. Tamo-Ye-Barg Codes

We analyze the I/O cost of a slightly modified construction in [10], whose repair bandwidth matches the cut-set bound of regenerating codes. Let $p_1 < p_2 < \dots < p_n$ be n prime numbers satisfying $r \mid (p_j - 1)$, $j \in [n]$, where $r \triangleq n - k$ (note that here we only consider the repair schemes where $d = n - 1$ helpers participate). Let α_j be an element of order p_j over \mathbb{F}_q , so that adjoining α_j to \mathbb{F}_q results in $\mathbb{F}_{q^{p_j}} = \mathbb{F}_q(\alpha_j)$, for every $j \in [n]$. Set $\mathcal{A} = \{\alpha_j\}_{j=1}^n$ to be the set of n (distinct) evaluation points of the code. Let β be an element of order r over $\mathbb{F}_{q^{p_1 p_2 \dots p_n}}$. The code $\text{RS}_{\mathbb{F}_{q^{\ell}}}(\mathcal{A}, k)$ is defined over $\mathbb{F}_{q^{\ell}} = \mathbb{F}_q(\alpha_1, \dots, \alpha_n, \beta)$ where $\ell = r \prod_{j=1}^n p_j$.

Assume that the failed node stores $c_{j^*} = f(\alpha_{j^*})$ for $f \in \mathbb{F}_{q^{\ell}}[x]$, $\deg(f) < k$. We consider the repair scheme \mathcal{R}_{j^*} based on a set of ℓ polynomials $\{\eta_{j^*, t} x^{w-1} : t \in [\ell/r], w \in [r]\}$, which are of degrees at most $r-1$. The coefficients $\eta_{j^*, t} \in \mathbb{F}_{q^{\ell}}$ are constructed as in [10] so that the set $\{\eta_{j^*, t} \alpha_{j^*}^{w-1} : t \in [\ell/r], w \in [r]\}$ forms an \mathbb{F}_q -basis of $\mathbb{F}_{q^{\ell}}$, which then qualifies \mathcal{R}_{j^*} as a repair scheme for c_{j^*} . Note that while in [10], the repair scheme for c_{j^*} is considered over the base field $\mathbb{F}_{q^{\prod_{j \neq j^*} p_j}}$, we go further down to \mathbb{F}_q instead and hence new factors will be introduced to reflect that change.

We define $\{\eta_{j^*, t}\}_{t=1}^{\ell/r}$ as the set

$$\left\{ \beta^u \alpha_{j^*}^{u+vr} \prod_{j \neq j^*} \alpha_j^{m_j} : 0 \leq u \leq r-1, 0 \leq v \leq \frac{p_{j^*} - 1}{r} - 1, 0 \leq m_j \leq p_j - 1 \right\} \cup \left\{ \left(\sum_{s=0}^{r-1} \beta^s \right) \alpha_{j^*}^{p_{j^*} - 1} \prod_{j \neq j^*} \alpha_j^{m_j} : 0 \leq m_j \leq p_j - 1 \right\}.$$

As the first and the second sets in the union have sizes $\frac{\ell}{r p_{j^*}} (p_{j^*} - 1)$ and $\frac{\ell}{r p_{j^*}}$, respectively, the union has size $\frac{\ell}{r}$ for a fixed j^* , as claimed. Compared to [10], the set $\{\prod_{j \neq j^*} \alpha_j^{m_j} : 0 \leq m_j \leq p_j - 1\}$ includes our newly introduced factors and serves as an \mathbb{F}_q -basis for $\mathbb{F}_{q^{\prod_{j \neq j^*} p_j}}$. The helper node storing $c_{j'} = f(\alpha_{j'})$, $j' \neq j^*$, chooses a basis $\mathcal{B}_{\alpha_{j'}}$ so that

$$\mathcal{B}'_{\alpha_{j'}} = \left\{ \beta^u \prod_{j=1}^n \alpha_j^{m_j} : 0 \leq u \leq r-1, 0 \leq m_j \leq p_j - 1 \right\}.$$

Similar to the two-coset construction, even if $\mathcal{B}'_{\alpha_{j'}}$ does not depend on $\alpha_{j'}$, the bases of the nodes can be different due to the different column multipliers. As discussed in Section IV-A, the I/O cost incurred at this node is equal to the number

of elements in $\mathcal{B}'_{\alpha_{j'}}$ required to generate $\{\eta_{j^*, t} \alpha_{j'}^{w-1} : t \in [\ell/r], w \in [r]\}$. Note that as $\alpha_{j'}$ is of order $p_{j'}$ over \mathbb{F}_q , all of its powers can be represented as an \mathbb{F}_q -linear combination of $\{\alpha_{j'}^{m_{j'}} : 0 \leq m_{j'} \leq p_{j'} - 1\}$. Therefore, every element in the set $\{\eta_{j^*, t} \alpha_{j'}^{w-1} : t \in [\ell/r], w \in [r]\}$ can be linearly generated over \mathbb{F}_q by the set

$$\left\{ \beta^u \alpha_{j^*}^{u+vr} \prod_{j \neq j^*} \alpha_j^{m_j} : 0 \leq u \leq r-1, 0 \leq v \leq \frac{p_{j^*} - 1}{r} - 1, 0 \leq m_j \leq p_j - 1 \right\} \cup \left\{ \beta^u \alpha_{j^*}^{p_{j^*} - 1} \prod_{j \neq j^*} \alpha_j^{m_j} : 0 \leq u \leq r-1, 0 \leq m_j \leq p_j - 1 \right\}.$$

Clearly, this is a subset of $\mathcal{B}'_{\alpha_{j'}}$ of cardinality $\frac{p_{j^*} + r - 1}{p_{j^*}} \frac{\ell}{r}$, which is the I/O cost incurred at the node storing $c_{j^*} = f(\alpha_{j^*})$. Therefore, the I/O cost in repairing $c_{j^*} = f(\alpha_{j^*})$ is $(1 + \frac{r-1}{p_{j^*}}) \frac{(n-1)\ell}{r}$ sub-symbols over \mathbb{F}_q , where $\frac{(n-1)\ell}{r}$ is the repair bandwidth of this repair scheme. Since $r \leq p_{j^*} - 1$ due to the assumption that $r \mid (p_{j^*} - 1)$, the I/O cost is at most twice the (optimal) repair bandwidth for any c_{j^*} , and the average I/O cost over the failed nodes is fairly close to the bandwidth.

REFERENCES

- [1] I. S. Reed and G. Solomon, "Polynomial codes over certain finite fields," *J. Soc. Ind. Appl. Math.*, vol. 8, no. 2, pp. 300–304, 1960.
- [2] H. Dau, I. Duursma, H. M. Kiah, and O. Milenkovic, "Repairing Reed-Solomon codes with multiple erasures," *IEEE Trans. Inform. Theory*, vol. 54, no. 10, pp. 6567–6582, 2018.
- [3] "HDFS Erasure Coding," <https://hadoop.apache.org/docs/r3.0.0/hadoop-project-dist/hadoop-hdfs/HDFSErasureCoding.html>.
- [4] K. Shanmugam, D. S. Papailiopoulos, A. G. Dimakis, and G. Caire, "A repair framework for scalar MDS codes," *IEEE J. Selected Areas Comm. (JSAC)*, vol. 32, no. 5, pp. 998–1007, 2014.
- [5] V. Guruswami and M. Wootters, "Repairing Reed-Solomon codes," in *Proc. Annu. Symp. Theory Comput. (STOC)*, 2016.
- [6] —, "Repairing Reed-Solomon codes," *IEEE Trans. Inform. Theory*, vol. 63, no. 9, pp. 5684–5698, 2017.
- [7] M. Ye and A. Barg, "Explicit constructions of MDS array codes and RS codes with optimal repair bandwidth," in *Proc. IEEE Int. Symp. Inform. Theory (ISIT)*, 2016, pp. 1202–1206.
- [8] H. Dau and O. Milenkovic, "Optimal repair schemes for some families of Reed-Solomon codes," in *Proc. IEEE Int. Symp. Inform. Theory (ISIT)*, 2017, pp. 346–350.
- [9] I. Duursma and H. Dau, "Low bandwidth repair of the RS(10,4) Reed-Solomon code," in *Proc. Inform. Theory Applicat. Workshop (ITA)*, 2017.
- [10] I. Tamo, M. Ye, and A. Barg, "Optimal repair of Reed-Solomon codes: Achieving the cut-set bound," in *Proc. 58th Annual IEEE Symp. Foundations Computer Sci. (FOCS)*, 2017.
- [11] —, "The repair problem for Reed-Solomon codes: Optimal repair of single and multiple erasures with almost optimal node size," in *IEEE Trans. Inform. Theory*, 2018, to appear.
- [12] H. Dau, I. Duursma, H. M. Kiah, and O. Milenkovic, "Repairing Reed-Solomon codes with two erasures," in *Proc. IEEE Int. Symp. Inform. Theory (ISIT)*, 2017, pp. 351–355.
- [13] B. Bartan and M. Wootters, "Repairing multiple failures for scalar MDS codes," in *Proc. 55th Annual Allerton Conf. Comm. Control Comput. (Allerton)*, 2017.
- [14] J. Mardia, B. Bartan, and M. Wootters, "Repairing multiple failures for scalar MDS codes," *IEEE Trans. Inform. Theory*, vol. 65, no. 5, pp. 2661–2672, 2018.
- [15] A. Chowdhury and A. Vardy, "Improved schemes for asymptotically optimal repair of MDS codes," in *Proc. 55th Annual Allerton Conf. Comm. Control Comput. (Allerton)*, 2017.
- [16] W. Li, Z. Wang, and H. Jafarkhani, "A tradeoff between the sub-packetization size and the repair bandwidth for Reed-Solomon code," in *Proc. 55th Annual Allerton Conf. Comm. Control Comput. (Allerton)*, 2017, pp. 942–949.
- [17] —, "A tradeoff between the sub-packetization size and the repair bandwidth for Reed-Solomon codes," *arXiv:1806.00496*, 2018.
- [18] H. Dau, I. Duursma, and H. Chu, "On the I/O costs of some repair schemes for full-length Reed-Solomon codes," in *Proc. IEEE Int. Symp. Inform. Theory (ISIT)*, 2018, pp. 1700–1704.
- [19] H. Dau and E. Viterbo, "Repair schemes with optimal I/O costs for full-length Reed-Solomon codes with two parities," in *Proc. IEEE Inform. Theory Workshop (ITW)*, 2018, pp. 590–594.
- [20] R. Lidl and H. Niederreiter, *Introduction to Finite Fields and Their Applications*. Cambridge University Press, 1986.
- [21] R. Roth, *Introduction to Coding Theory*. New York, NY, USA: Cambridge University Press, 2006.