# On Optimal Locally Repairable Codes With Multiple Disjoint Repair Sets

Han Cai, *Member, IEEE*, Ying Miao, Moshe Schwartz, *Senior Member, IEEE*,
and Xiaohu Tang, *Senior Member, IEEE*

*Abstract*—Locally repairable codes are desirable for distributed storage systems to improve the repair efficiency. In this paper, a new combination of codes with locality and codes with multiple disjoint repair sets (also called availability) is introduced. Accordingly, a Singleton-type bound is derived for the new code, which contains those bounds in [9], [20], [28] as special cases. Optimal constructions are proposed with respect to the new bound. In addition, these constructions can also generate optimal codes with multiple disjoint repair sets with respect to the bound in [28], which to the best of our knowledge, are the first explicit constructions that can achieve the bound in [28].

*Index Terms*—Availability, distributed storage, locally repairable code.

## I. INTRODUCTION

**N**OWADAYS, large-scale cloud storage and distributed file systems such as Amazon Elastic Block Store (EBS) and Google File System (GoogleFS) have reached such a massive scale that disk failures are the norm rather than the exception. One of the simplest solutions to protect data from disk failures in these systems is straightforward replication of data packets across different disks. However, this solution suffers from a larger storage overhead. To reduce the storage overhead, an alternative solution based on storage codes was proposed. An $[n, k]$ storage code encodes $k$ information symbols to $n$ symbols and stores them across $n$ disks in a storage system. Generally speaking, among all storage codes, maximum distance separable (MDS) codes are preferred for practical systems since they can lead to dramatic improvements, both in terms of redundancy and in terms of reliability, compared with replication [9]. Nevertheless, $[n, k]$ MDS codes have a drawback that whenever one wants to recover a symbol, one needs to contact $k$ surviving symbols, which is costly, especially in large-scale distributed file systems.

To overcome the above drawback, locally repairable codes were introduced to reduce the number of symbols contacted during the repair process. More precisely, the concept of locality for a code $\mathcal{C}$ was initially studied in [11] to ensure that a failed symbol can be recovered by only accessing other $r \ll k$ symbols which form a repair set [2].

However, the original concept of locality only works when exactly one erasure occurs. To guarantee that the system can locally recover from multiple erasures (say, $\delta - 1 > 1$ erasures), there are two main extensions in the literature. The first approach is to let the repair set contain $\delta - 1$ redundancies. In this case, even if $\delta - 1$ erasures occur, the failed symbols may still be recovered locally by the remaining symbols in the corresponding repair sets [20]. The second approach is to provide the code symbols with $\delta - 1$ disjoint repair sets [28]. In this scheme, if there are at most $\delta - 1$ erasures, then for each failed symbol at least one complete repair set can be accessed to recover the failed symbol locally. In particular, a code with multiple repair sets (also called availability [22]) has the advantage of good parallel reading ability, since each repair set can be seen as a backup for the target symbol and thus can be accessed independently.

Up to now, some upper bounds on the minimum Hamming distance of locally repairable codes have been derived, such as the Singleton-type bounds in [9], [19], [20], [29], bounds depending on the size of the alphabet [1], [3], the bound for locally repairable codes with multiple erasure tolerance [22], [28], etc. Numerous constructions of optimal locally repairable codes with respect to those bounds have been reported in the literature, e.g., see [2], [5], [8], [9], [11], [18]–[22], [24], [26], [28], and the references therein. All these bounds and constructions are either under the definition of locality in [20] or the one in [28] and [22].

In this paper, to allow the system to recover *locally* from multiple erasures, we go beyond the aforementioned solutions and establish a more general framework for locally repairable codes with multiple disjoint repair sets. Firstly, we combine the solutions in [20], [22], [28] by a trade-off between the number of repair sets and the number of redundancies in each repair set. As a result, the locally repairable codes in [20], [22], [28] are exactly the extremal cases of our setting. Secondly, we derive a new Singleton-type bound for

the generalized locally repairable codes, which contains the bounds in [9], [20], [28] as special cases. Finally, we describe constructions of optimal locally repairable codes with respect to the bound we derived (Corollaries 6 and 8). As a byproduct, the constructions can generate optimal locally repairable codes with multiple disjoint repair sets with respect to the bound in [28] (Corollaries 5 and 7). To the best of our knowledge, no explicit construction has achieved the bound in [28] before. As a comparison, we list the known optimal locally repairable codes with multiple repair sets in Table I.

The remainder of this paper is organized as follows. Section II introduces some preliminaries about locally repairable codes. Section III proposes a new definition for locality that generalizes those in [20] and [22], [28]. Section IV establishes a Singleton-type bound for locally repairable codes. Sections V and VI present constructions of optimal locally repairable codes with respect to the new bound. Section VII concludes this paper with some remarks.

## II. PRELIMINARIES

In this section we describe the notation used, and give a short overview of locally repairable codes. Throughout this paper, the following notation is used unless otherwise stated: If $n$ is a positive integer then $[n]$ denotes the set $\{1, 2, \cdots, n\}$. For integers $a > 0$ and $b$, $\langle b \rangle_a$ denotes the least nonnegative residue of $b$ modulo $a$.

We let $\mathbb{F}_q$ denote the finite field with $q$ elements, where $q$ is a prime power. An $[n, k]_q$ linear code $\mathcal{C}$ over $\mathbb{F}_q$ is a $k$-dimensional subspace of $\mathbb{F}_q^n$ with a $k \times n$ generator matrix $G = (\mathbf{g}_1, \mathbf{g}_2, \cdots, \mathbf{g}_n)$, where $\mathbf{g}_i$ is a column vector of dimension $k$ for all $1 \leqslant i \leqslant n$. We also call $\mathcal{C}$ an $[n, k, d]_q$ linear code when the minimal Hamming distance $d$ is available. Note that throughout this paper we only consider the Hamming distance. For a subset $S \subseteq [n]$, we use $|S|$, Span($S$), and Rank($S$) to denote the cardinality of $S$, the linear space spanned by $\{\mathbf{g}_i : i \in S\}$ over $\mathbb{F}_q$, and the dimension of Span($S$), respectively.

In [11], Huang *et al.* first studied the locality of code symbols via the Pyramid code. The $j$th ($1 \leqslant j \leqslant n$) code symbol, in an $[n, k]_q$ linear code $\mathcal{C}$, is said to have locality $r$ ($1 \leqslant r \leqslant k$), if it can be recovered by accessing at most $r$ other symbols in $\mathcal{C}$. More precisely:

**Definition 1** ([9])**:** For any column $\mathbf{g}_j$, $1 \leqslant j \leqslant n$, of a generator matrix $G$ of an $[n, k]_q$ linear code $\mathcal{C}$, define Loc($\mathbf{g}_j$) as the smallest integer $r$ such that there exists a set $R = \{j_1, j_2, \cdots, j_r\} \subseteq [n] \setminus \{j\}$ satisfying $\mathbf{g}_j \in$ Span($R$), i.e., there exist $\lambda_t \in \mathbb{F}_q$, $1 \leqslant t \leqslant r$ such that

$$\mathbf{g}_j = \sum_{t=1}^{r} \lambda_t \mathbf{g}_{j_t}. \qquad (1)$$

Define Loc($S$) = $\max_{j \in S}$ Loc($\mathbf{g}_j$) for any set $S \subseteq [n]$. The code $\mathcal{C}$ is said to have information locality $r$, if there exists $S \subseteq [n]$ with Rank($S$) = $k$ and Loc($S$) = $r$.

Obviously, $c_j = \sum_{t=1}^{r} \lambda_t c_{j_t}$ for every codeword $(c_1, c_2, \cdots, c_n) \in \mathcal{C}$ is equivalent with $\mathbf{g}_j = \sum_{t=1}^{r} \lambda_t \mathbf{g}_{j_t}$, where $\lambda_t \in \mathbb{F}_q$ for $1 \leqslant t \leqslant r$. Therefore, throughout this paper we do not distinguish between the $j$th code symbol (i.e., $c_j$ for any codeword $(c_1, c_2, \cdots, c_n) \in \mathcal{C}$ and the $j$th column of $\mathbf{g}_j$ of a generator matrix $G$ for $\mathcal{C}$. Thus, we call both $c_j$ and $\mathbf{g}_j$ as the $j$th code symbol for $1 \leqslant j \leqslant n$.

According to (1), a single erasure can be recovered by accessing at most other $r$ symbols. Two methods appear in the literature to generalize this by guaranteeing local recovery from more than one erasure. The first method is to let the repair set contain more than one redundancy, say $\delta - 1 > 1$ redundancies:

**Definition 2** ([20])**:** The $j$th column $\mathbf{g}_j$, $1 \leqslant j \leqslant n$, of the generator matrix $G$ of an $[n, k]_q$ linear code $\mathcal{C}$ is said to have $(r, \delta)$-locality, if there exists a subset $S_j \subseteq [n]$ such that:

- $j \in S_j$ and $|S_j| \leqslant r + \delta - 1$; and
- the minimum Hamming distance of the punctured code $\mathcal{C}|_{S_j}$ obtained by deleting the code symbols $c_t$ ($t \in [n] \setminus S_j$) is at least $\delta$,

where the set $S_j \setminus \{j\}$ is also called a repair set of $\mathbf{g}_j$. Further, the code $\mathcal{C}$ is said to have information $(r, \delta)$-locality if there exists $S \subseteq [n]$ with Rank($S$) = $k$ such that for each $j \in S$, $\mathbf{g}_j$ has $(r, \delta)$-locality.

In [20] the following upper bound on the minimum Hamming distance of linear codes with information $(r, \delta)$-locality was derived.

**Lemma 1** ( [20])**:** For an $[n, k, d]_q$ linear code with information $(r, \delta)$-locality,

$$d \leqslant n - k + 1 - \left( \left\lceil \frac{k}{r} \right\rceil - 1 \right)(\delta - 1). \qquad (2)$$

The second method to guarantee local recovery from multiple erasures is to provide code symbols with multiple pairwise disjoint repair sets, say $\delta - 1$ sets, of size at most $r$ [28], which are also called $(r, \delta)$-availability [22].

**Definition 3** ( [28], [22])**:** The $j$th column $\mathbf{g}_j$, $1 \leqslant j \leqslant n$, of a generator matrix of an $[n, k, d]_q$ linear code $\mathcal{C}$ is said to have $(r, \delta)_c$-locality, or $(r, \delta)$-availability, if there exist $\delta - 1$ pairwise disjoint sets $R_1^j, R_2^j, \ldots, R_{\delta-1}^j \subseteq [n] \setminus \{j\}$, satisfying

- $\left| R_t^j \right| \leqslant r$; and
- $\mathbf{g}_j \in$ Span $\left( R_t^j \right)$

for all $1 \leqslant t \leqslant \delta - 1$, where each $R_t^j$ is called a repair set of $\mathbf{g}_j$. Furthermore, the code $\mathcal{C}$ is said to have information $(r, \delta)_c$-locality if there is a subset $S \subseteq [n]$ with Rank($S$) = $k$ such that $\mathbf{g}_j$ has $(r, \delta)_c$-locality for each $j \in S$.

In this scheme, if there are at most $\delta - 1$ erasures, then for each erased symbol at least one complete repair set can be accessed to recover it locally. Each repair set $R_t^j$ can be viewed as a backup for the target code symbol $\mathbf{g}_j$ and hence these pairwise disjoint repair sets can be accessed independently, which means that $\mathbf{g}_j$ has parallel reading ability. The minimum Hamming distance $d$ of a linear code $\mathcal{C}$ with information $(r, \delta)_c$-locality is upper bounded as follows.

**Lemma 2** ([28])**:** For an $[n, k, d]_q$ linear code with information $(r, \delta)_c$-locality,

$$d \leqslant n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil.$$

TABLE I
KNOWN OPTIMAL LOCALLY REPAIRABLE CODES WITH MULTIPLE REPAIR SETS

| Parameters | Locality | Field size | Constraints | Explicit Construction | Ref. |
|---|---|---|---|---|---|
| $[n,k,d]_q$ | $(r,\delta)_c$ | $q > 1 + \binom{n}{k+\sigma}$ | $\sigma = \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$ $n \geqslant k(1+r(\delta-1))$ | No | [28] |
| $[n,k,d]_q$ | $(r,\delta)_c$ | $q \geqslant q_1^{k(1+(r-1)(\delta-1))}$ | $q_1 \geqslant r+1$ $n = k(1+r(\delta-1))$ | Yes | Construction A (Corollary 5) |
| $[n,k,d]_q$ | $(r,\delta)_c$ | $q \geqslant q_1^{1+(r-1)N}$ | $q_1 \geqslant k+1$ $n = k(1+r(\delta-1))$ | Yes | Construction B (Corollary 7) |
| $[n,k,d]_q$ | $(r,N,\delta)$ | $q \geqslant q_1^{k(1+(r-1)N)}$ | $\delta-1 = N(d^*-1),\ d^* > 2$ $q_1 \geqslant r+1,\ n = k(1+r(\delta-1))$ | Yes | Construction A (Corollary 6) |
| $[n,k,d]_q$ | $(r,N,\delta)$ | $q \geqslant q_1^{1+(r-1)N}$ | $\delta-1 = N(d^*-1),\ n = k(1+r(\delta-1))$ $q_1 \geqslant \max\{r+d^*-1, k+1\},\ d^* > 2$ | Yes | Construction B (Corollary 8) |

We conclude this section with three remarks concerning the two definitions for locality and their connection to previous literature.

**Remark 1:** When $\delta = 2$, both definitions of $(r,\delta)$-locality and $(r,\delta)_c$-locality coincide with Definition 1 from [11]. Both codes with information $(r,\delta)$-locality and codes with information $(r,\delta)_c$-locality with $\delta > 2$ can recover an information symbol with the help of at most $r$ surviving symbols when there are at most $\delta - 1$ erasures [20], [28].

**Remark 2:** In [28], Wang and Zhang proved that the bound in Lemma 2 can be achieved when the code rate is low and the underlying finite field is sufficiently large. Later, in [27], Tamo *et al.* derived a new bound for codes with $(r,\delta)_c$-locality, which improves the bound in Lemma 2 for the high code-rate case.

**Remark 3:** Optimal constructions for locally repairable codes with respect to the bound in Lemma 1 may be found, for example, in [9], [21], [24], [26]. Compared with the $(r,\delta)$-locality, codes with information $(r,\delta)_c$-locality have the advantage of good parallel reading ability [22]. However, to the best of our knowledge, no explicit construction achieves the bound in Lemma 2. One severely limited solution for locally repairable codes with $(r,\delta)_c$-locality assumes that each repair set contains exactly one check symbol. For a bound and corresponding optimal constructions for this limited setting the reader is referred to [4], [10], [18], [22], [25].

## III. A GENERAL DEFINITION FOR LOCALLY REPAIRABLE CODES

We give a definition for locality which generalizes previous definitions, and prove that it indeed guarantees local recovery from multiple erasures. By Definitions 2 and 3, the $(r,\delta)$- or $(r,\delta)_c$-locality properties both guarantee local recovery from most $\delta - 1$ erasures. However, they provide different availability for code symbols and the trade-off between the parameters are also different by Lemmas 1 and 2. The motivation for this study is to find the trade-off between availability and the repair ability for each repair set when the code symbols can be locally recovered from $\delta - 1$ erasures. To this end, we first generalize the definition for symbol locality that can guarantee local recovery from $\delta - 1$ erasures.

**Definition 4:** The $j$th column $\mathbf{g}_j$, $1 \leqslant j \leqslant n$, of a generator matrix of an $[n,k,d]_q$ linear code $\mathcal{C}$, is said to have $(r, N_j, \delta)$-*locality*, if there exist $N_j \geqslant 1$ pairwise disjoint repair sets, i.e., $N_j \geqslant 1$ pairwise disjoint subsets of $\{\mathbf{g}_i : 1 \leqslant i \leqslant n\} \setminus \{\mathbf{g}_j\}$, denoted $R_1^j, R_2^j, \cdots, R_{N_j}^j$, that satisfy the following conditions:

- For any $1 \leqslant l \leqslant N_j$, $\left| R_l^j \right| \leqslant r + d_l^j - 2$;
- For any $1 \leqslant l \leqslant N_j$, the code $\mathcal{C}|_{R_l^j \cup \{\mathbf{g}_j\}}$ is a linear code with minimum Hamming distance $d_l^j \geqslant 2$;
- $\sum_{1 \leqslant l \leqslant N_j} \left( d_l^j - 1 \right) \geqslant \delta - 1$.

Furthermore, the code $\mathcal{C}$ is said to have information $(r, \mathbf{N}, \delta)$-locality, if there is a subset $S = \{s_1, s_2, \ldots, s_k\} \subseteq [n]$ with $1 \leqslant s_1 < s_2 < \cdots < s_k \leqslant n$ and $\text{Rank}(S) = k$ such that $\mathbf{g}_j$ has $(r, N_j, \delta)$-locality for each $j \in S$, where $\mathbf{N} = (N_{s_1}, N_{s_2}, \ldots, N_{s_k})$.

**Remark 4:** The first two conditions for the $(r, N_j, \delta)$-locality are used to make sure that each $R_i^j$ for $1 \leqslant i \leqslant N_j$ is capable of recovering $\mathbf{g}_j$ by only accessing $r$ symbols. The first two conditions also mean that $\mathbf{g}_j$ has availability $N_j$, i.e., allowing $N_j + 1$ parallel reads for the code symbol $\mathbf{g}_j$, since each repair set can be read in parallel to recover $\mathbf{g}_j$. The last restriction guarantees the recovery from $\delta - 1$ erasures. The symbol $\mathbf{g}_j$ can be recovered after $\delta - 1$ erasures since regardless of the way those erasures are distributed over $N_j$ pairwise disjoint repair sets, at least one repair set say the $l$-th, is not hit by more than $d_l^j$ erasures. Thus, we can recover $\mathbf{g}_j$ locally. Refer to Lemma 3 and its proof for more details.

Based on the above definition, we fix the following notation for an $[n,k,d]_q$ code with $(r, \mathbf{N}, \delta)$-locality throughout this paper:

- $\mathbf{I}_j$ denotes the $j$th information symbol for $1 \leqslant j \leqslant k$. Without loss of generality, we assume that they are exactly the first $k$ code symbols, that is, $\mathbf{I}_j = \mathbf{g}_j$ for $1 \leqslant j \leqslant k$;
- $R_1^j, R_2^j, \cdots, R_{N_j}^j$ denote the $N_j$ pairwise disjoint repair sets for $\mathbf{I}_j$, $1 \leqslant j \leqslant k$;
- $U_j$ denotes the union of $\{\mathbf{I}_j\}$ and all pairwise disjoint repair sets for $\mathbf{I}_j$, i.e.,

$$U_j = \{\mathbf{I}_j\} \cup \left( \bigcup_{1 \leqslant l \leqslant N_j} R_l^j \right) \text{ for } 1 \leqslant j \leqslant k. \quad (3)$$

Accordingly, we say the code $\mathcal{C}$ has information $(r, \mathbf{N}, \delta)$-locality, if $\mathbf{I}_j$ has $(r, N_j, \delta)$-locality for each $1 \leqslant j \leqslant k$, where

$\mathbf{N} = (N_1, N_2, \ldots, N_k)$. When $\mathbf{N} = (a, a, \ldots, a)$, we denote it as information $(r, a, \delta)$-locality.

**Lemma 3:** Let $\mathcal{C}$ be a linear code with information $(r, \mathbf{N}, \delta)$-locality, and let $E$ be an erasure pattern. If $|E| \leqslant \delta - 1$, then the information symbols in $E$ can be recovered locally, i.e., recovered by accessing at most $r$ surviving symbols.

*Proof:* We assume to the contrary that there exists an erased information symbol, say $\mathbf{I}_j \in E$, which cannot be recovered locally. Then for $1 \leqslant l \leqslant N_j$, $|E \cap (R_l^j \cup \{\mathbf{I}_j\})| \geqslant d_l^j$, otherwise the symbols in $R_l^j \setminus E$ can be accessed to recover $\mathbf{I}_j$ locally since the code $C|_{R_l^j \cup \{\mathbf{I}_j\}}$ is a linear code with minimum Hamming distance $d_l^j$ in Definition 4. Now the fact $\mathbf{I}_j \notin R_l^j$ means $|(E \setminus \{\mathbf{I}_j\}) \cap R_l^j| \geqslant d_l^j - 1$ for $1 \leqslant l \leqslant N_j$. Thus,

$$|E| = 1 + |E \setminus \{\mathbf{I}_j\}| \geqslant 1 + \left| \bigcup_{1 \leqslant l \leqslant N_j} \left( E \cap R_l^j \right) \setminus \{\mathbf{I}_j\} \right|$$
$$\geqslant 1 + \sum_{1 \leqslant l \leqslant N_j} \left( d_l^j - 1 \right) \geqslant \delta,$$

where the last inequality holds by Definition 4, a contradiction. ∎

In particular, it is easily seen from Definition 4 that:
- The $(r, 1, 2)$-locality in Definition 4 corresponds to the $r$-locality in Definition 1.
- The $(r, 1, \delta)$-locality in Definition 4 corresponds to the $(r, \delta)$-locality in Definition 2.
- The $(r, \delta - 1, \delta)$-locality for the case $d_l^j = 2$, $1 \leqslant j \leqslant \delta - 1$, $1 \leqslant l \leqslant k$ in Definition 4 corresponds to the $(r, \delta)_c$-locality in Definition 3.

In summary, the definitions in [9], [20] and [28] correspond to two extremal cases of Definition 4. For any given $r$ and $\delta$, a code $\mathcal{C}$ with information $(r, \mathbf{N}, \delta)$-locality can locally repair a failed information symbol by accessing at most $r$ other symbols when at most $\delta - 1$ erasures occur. Specifically, different $N_j$ means different numbers of repair sets for $\mathbf{I}_j$, $1 \leqslant j \leqslant k$, i.e., different parallel reading abilities. Thus, Definition 4 not only contains the two previous definitions for locality as special cases, but also suggests the existence of new scenarios in which local recovery is possible. As a comparison, in Figure 1, we draw an illustration of different types of localities with the property that $\mathbf{g}_j$ can be recovered by 4 symbols when there are at most 5 erasures.

## IV. THE BOUND FOR LINEAR CODES WITH INFORMATION $(r, \mathbf{N}, \delta)$-LOCALITY

The goal of this section is to establish an upper bound on the minimum Hamming distance of linear codes with information $(r, \mathbf{N}, \delta)$-locality. This bound appears in Theorem 1. In order to prove the bound, a careful analysis of subsets of codeword coordinates is performed in Lemma 4 and Lemma 5, tying together the size of subsets of coordinates and their rank. Following the main result of this section, several corollaries are given, studying various specific sets of parameters implied by the result of Theorem 1.

We start with folklore and known results:

**Fact 1:** Let $W$ and $S$ be two sets of vectors over $\mathbb{F}_q$ with $S \subseteq W$. Then, $|W| - |S| \geqslant \text{Rank}(W) - \text{Rank}(S)$.

**Lemma 4** ( [15]): An $[n, k]_q$ linear code $\mathcal{C}$ has a minimum Hamming distance $d$ if and only if $d$ is the largest integer such that

$$|S| \leqslant n - d$$

for every $S \subseteq \{\mathbf{g}_j : j \in [n]\}$ with $\text{Rank}(S) \leqslant k - 1$.

In addition, the following results will be used frequently in proving our bound.

**Lemma 5:** Let $\mathcal{C}$ be an $[n, k]_q$ linear code with information $(r, \mathbf{N}, \delta)$-locality, and let $U_j$ be defined by (3) for $1 \leqslant j \leqslant k$.

1) If $S \subseteq U_j$ $(1 \leqslant j \leqslant k)$, then

$$\text{Rank}(S) \leqslant 1 + \sum_{|S \cap R_l^j| \geqslant r} (r - 1) + \sum_{|S \cap R_l^j| < r} \left| S \cap R_l^j \right|; \quad (4)$$

2) If $S \subseteq \{\mathbf{g}_j : j \in [n]\}$ and

$$\left| R_l^j \cap S \right| \leqslant \left| R_l^j \right| - d_l^j + 1 \quad (5)$$

for $1 \leqslant l \leqslant \Lambda \leqslant N_j$, then

$$\left| \left( \bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j \right) \cup S \right| - |S|$$
$$\geqslant \text{Rank}\left( \left( \bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j \right) \cup S \right)$$
$$- \text{Rank}(S) + \sum_{1 \leqslant l \leqslant \Lambda} \left( d_l^j - 1 \right). \quad (6)$$

Particularly, if $\Lambda = N_j$ then

$$\left| U_j \cup S \right| - |S| \geqslant \text{Rank}\left( U_j \cup S \right) - \text{Rank}(S) + \delta - 1. \quad (7)$$

*Proof:* First, we state the following property:

P1. Any set $R_l^j \cup \{\mathbf{I}_j\}$ can be spanned by any of their $r$ symbols, in particular, $\mathbf{I}_j$ and any other $r - 1$ symbols from $R_l^j$.

Property P1 holds since $\mathcal{C}|_{R_l^j \cup \{\mathbf{I}_j\}}$ is a linear code with minimum Hamming distance $d_l^j$ and $|R_l^j \cup \{\mathbf{I}_j\}| \leqslant r + d_l^j - 1$ for any $1 \leqslant l \leqslant N_j$.

Thus, for the first part, according to (3) and P1,

$$\text{Rank}(S) \leqslant \text{Rank}(\{\mathbf{I}_j\}) + \sum_{|S \cap R_l^j| \geqslant r} \left( \text{Rank}\left( R_l^j \right) - 1 \right)$$
$$+ \sum_{|S \cap R_l^j| < r} \text{Rank}\left( S \cap R_l^j \right)$$
$$\leqslant 1 + \sum_{|S \cap R_l^j| \geqslant r} (r - 1) + \sum_{|S \cap R_l^j| < r} \left| S \cap R_l^j \right|.$$

Fig. 1. A comparison between different types of localities for the $j$th code symbol $\mathbf{g}_j$, where the curves and lines drawn with dashed lines correspond to repair sets of $\mathbf{g}_j$ which satisfy that any 4 points suffice to recover the curve.

For the second part, (5) and P1 then mean that we can find $D_l^j \subseteq (R_l^j \setminus S)$ with $|D_l^j| = d_l^j - 1$ such that

$$\text{Rank}\left(\left(R_l^j \cup \{\mathbf{I}_j\}\right) \setminus D_l^j\right)$$
$$= \text{Rank}\left(R_l^j \cup \{\mathbf{I}_j\}\right), \ 1 \leqslant l \leqslant N_j. \qquad (8)$$

Note from Definition 4 that $R_{l_1}^j \cap R_{l_2}^j = \emptyset$ for $1 \leqslant l_1 < l_2 \leqslant N_j$, thus $D_{l_1}^j \cap D_{l_2}^j = \emptyset$ for $1 \leqslant l_1 < l_2 \leqslant N_j$, i.e.,

$$\left| \bigcup_{1 \leqslant l \leqslant \Lambda} D_l^j \right| = \sum_{1 \leqslant l \leqslant \Lambda} \left| D_l^j \right| = \sum_{1 \leqslant l \leqslant \Lambda} \left( d_l^j - 1 \right). \qquad (9)$$

Set

$$W = \left(\left(\bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j\right) \cup S\right) \setminus \left(\bigcup_{1 \leqslant l \leqslant N_j} D_l^j\right).$$

It follows from (8) and (9) that

$$|W| = \left|\left(\bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j\right) \cup S\right| - \sum_{1 \leqslant l \leqslant \Lambda} \left(d_l^j - 1\right)$$

and

$$\text{Rank}(W) = \text{Rank}\left(\left(\bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j\right) \cup S\right).$$

Thus, applying Fact 1 to $S \subset W$, we have

$$\left|\left(\bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j\right) \cup S\right| - |S| = |W| - |S| + \sum_{1 \leqslant l \leqslant \Lambda} \left(d_l^j - 1\right)$$

$$\geqslant \text{Rank}\left(\left(\bigcup_{1 \leqslant l \leqslant \Lambda} R_l^j\right) \cup S\right)$$

$$- \text{Rank}(S) + \sum_{1 \leqslant l \leqslant \Lambda} \left(d_l^j - 1\right),$$

which turns out to be (7) when $\Lambda = N_j$ because of (3) and $\sum_{1 \leqslant l \leqslant N_j} \left(d_l^j - 1\right) \geqslant \delta - 1$. $\blacksquare$

Now, we are ready to present our bound.

**Theorem 1:** For any $[n, k, d]_q$ linear code with information $(r, \mathbf{N}, \delta)$-locality,

$$d \leqslant \begin{cases} n - k + 1 - \mu(\delta - 1), \\ \qquad\qquad \text{if } (1 + N(r-1))|(k-1), \\ n - k + 1 - \mu(\delta - 1) - \left\lceil \frac{\Lambda(\delta-1)}{N} \right\rceil, \ \text{otherwise,} \end{cases}$$
$$\qquad (10)$$

where $N = \max(\{N_j : 1 \leqslant j \leqslant k\})$, $\mu = \lfloor \frac{k-1}{1+N(r-1)} \rfloor$, and $\Lambda = \lfloor \frac{\langle k-1 \rangle_{1+N(r-1)} - 1}{r-1} \rfloor$.

*Proof:* According to Lemma 4, to prove this theorem it suffices to find a set $S$ with rank $k-1$ and

$$|S| \geqslant \begin{cases} k - 1 + \mu(\delta - 1), \\ \qquad\qquad \text{if } (1 + N(r-1))|(k-1), \qquad (11) \\ k - 1 + \mu(\delta - 1) + \left\lceil \frac{\Lambda(\delta-1)}{N} \right\rceil, \quad \text{otherwise.} \end{cases}$$

If $\mu = \lfloor \frac{k-1}{1+N(r-1)} \rfloor = 0$, set $S_0 = \emptyset$. Otherwise, we can select $\mu$ information symbols, say $\mathbf{I}_{j_1}, \mathbf{I}_{j_2}, \cdots, \mathbf{I}_{j_\mu}$, such that $\mathbf{I}_{j_i} \notin \text{Span}(S_{i-1})$, where $S_0 = \emptyset$ and $S_i = \bigcup_{1 \leqslant l \leqslant i} U_{j_l}$ for $1 \leqslant i \leqslant \mu$. This is to say, $S_{i+1} = S_i \cup U_{j_{i+1}}$ for $0 \leqslant i < \mu$. Then, for $0 \leqslant i < \mu$, $\text{Rank}(S_{i+1}) \geqslant \text{Rank}(S_i) + 1$ and

$$\text{Rank}(S_\mu) \leqslant \sum_{l=1}^{\mu} \text{Rank}(U_{j_l})$$
$$\leqslant \mu(1 + N(r-1)) \leqslant k - 1, \qquad (12)$$

where we use the inequality $\text{Rank}(U_j) \leqslant 1 + N_j(r-1) \leqslant 1 + N(r-1)$ by (4).

Recall from (12) that $\text{Rank}(S_\mu) = k - 1$ only if $(1 + N(r-1)) \mid (k-1)$. Thus, if $(1+N(r-1)) \nmid (k-1)$, there is one more information symbol $\mathbf{I}_{j_{\mu+1}}$ such that $\mathbf{I}_{j_{\mu+1}} \notin \text{Span}(S_\mu)$. When $\Lambda \geqslant 1$ and $N_{j_{\mu+1}} \geqslant \Lambda$, among $\binom{N_{j_{\mu+1}}}{\Lambda}$ distinct $\Lambda$-sets, each $R_l^{j_{\mu+1}}$ ($1 \leqslant l \leqslant N_{j_{\mu+1}}$) appears $\binom{N_{j_{\mu+1}}-1}{\Lambda-1}$ times.

According to the pigeonhole principle, there must exist $\Lambda$ repair sets, say $R_l^{j_{\mu+1}}$ for $1 \leqslant l \leqslant \Lambda$, such that

$$\sum_{1 \leqslant l \leqslant \Lambda} \left( d_l^{j_{\mu+1}} - 1 \right) \geqslant \left\lceil \frac{\binom{N_{j_{\mu+1}}-1}{\Lambda-1}(\delta-1)}{\binom{N_{j_{\mu+1}}}{\Lambda}} \right\rceil$$

$$= \left\lceil \frac{\Lambda(\delta-1)}{N_{j_{\mu+1}}} \right\rceil. \tag{13}$$

In this case, i.e., $(1 + N(r-1)) \nmid (k-1)$ and $\Lambda \geqslant 1$, set

$$S_{\mu+1} = S_\mu \cup \{\mathbf{I}_{j_{\mu+1}}\} \cup \left( \bigcup_{1 \leqslant l \leqslant \min\{\Lambda, N_{j_{\mu+1}}\}} R_l^{j_{\mu+1}} \right).$$

Then, we have

$$\begin{aligned} &\text{Rank}(S_{\mu+1}) \\ &\leqslant \text{Rank}(S_\mu) \\ &\quad + \text{Rank}\left( \mathbf{I}_{j_{\mu+1}} \cup \left( \bigcup_{1 \leqslant l \leqslant \min\{\Lambda, N_{j_{\mu+1}}\}} R_l^{j_{\mu+1}} \right) \right) \\ &\leqslant \mu(1 + N(r-1)) + 1 + \Lambda(r-1) \\ &\leqslant k-1, \end{aligned}$$

where we use (4) and (12) in the second inequality and the fact $\Lambda = \lfloor \frac{(k-1)_{1+N(r-1)}-1}{r-1} \rfloor$ in the third inequality, respectively.

Note that $\mathbf{I}_{j_i} \notin \text{Span}(S_{i-1})$, which implies that (5) holds for all $S_{i-1}$ and $R_l^{j_i}$, $1 \leqslant i \leqslant \mu+1$ and $1 \leqslant l \leqslant N_{j_i}$. Otherwise, the fact that $\mathcal{C}|_{R_l^{j_i} \cup \{\mathbf{I}_{j_i}\}}$ has minimum Hamming distance $d_l^{j_i}$ leads to $\mathbf{I}_{j_i} \in \text{Span}(S_{i-1})$, a contradiction. Therefore, applying (7) in place of $S = S_0, \cdots, S_\mu$ sequentially, we have

$$\begin{aligned} |S_\mu| &= \sum_{i=0}^{\mu-1} (|S_{i+1}| - |S_i|) \\ &\geqslant \sum_{i=0}^{\mu-1} (\text{Rank}(S_{i+1}) - \text{Rank}(S_i)) + \mu(\delta-1) \\ &\geqslant \text{Rank}(S_\mu) + \mu(\delta-1), \end{aligned} \tag{14}$$

where we use the fact that $|S_0| = \text{Rank}(S_0) = 0$ due to $S_0 = \emptyset$. Moreover, when $(1 + N(r-1)) \nmid (k-1)$, by applying (6) we can get

$$\begin{aligned} |S_{\mu+1}| - |S_\mu| &\geqslant \text{Rank}(S_{\mu+1}) - \text{Rank}(S_\mu) \\ &\quad + \sum_{1 \leqslant l \leqslant \min\{\Lambda, N_{j_{\mu+1}}\}} \left( d_l^{j_{\mu+1}} - 1 \right) \\ &\geqslant \text{Rank}(S_{\mu+1}) - \text{Rank}(S_\mu) \\ &\quad + \begin{cases} \left\lceil \frac{\Lambda(\delta-1)}{N_{j_{\mu+1}}} \right\rceil, & \text{if } N_{j_{\mu+1}} \geqslant \Lambda \\ \delta-1, & \text{if } N_{\mu+1} < \Lambda \end{cases} \\ &\geqslant \text{Rank}(S_{\mu+1}) - \text{Rank}(S_\mu) + \left\lceil \frac{\Lambda(\delta-1)}{N} \right\rceil, \end{aligned}$$

where, for $\Lambda \geqslant 1$, we use (13) for the case $N_{j_{\mu+1}} \geqslant \Lambda$ and $\sum_{1 \leqslant l \leqslant N_{j_{\mu+1}}} (d_l^{j_{\mu+1}} - 1) \geqslant \delta - 1$ for the case $N_{j_{\mu+1}} < \Lambda$. Then,

together with (14) gives

$$|S_{\mu+1}| \geqslant \text{Rank}(S_{\mu+1}) + \mu(\delta-1) + \left\lceil \frac{\Lambda(\delta-1)}{N} \right\rceil.$$

Finally, form a set $S$ with $\text{Rank}(S) = k-1$ by appending some elements into $S_\mu$ if $(1 + N(r-1)) \mid (k-1)$ or $\Lambda = 0$, and $S_{\mu+1}$ otherwise. Then, the desired result (11) follows from Fact 1. ∎

When $N = 1$, Theorem 1 is exactly the bound in Lemma 1, first derived in [20] ( [9] for $\delta = 2$).

**Corollary 1:** For any $[n, k, d]_q$ linear code with information $(r, N = 1, \delta)$-locality,

$$d \leqslant n-k+1 - \left( \left\lceil \frac{k}{r} \right\rceil - 1 \right)(\delta-1). \tag{15}$$

*Proof:* For the case $N = 1$, it is easy to see that $\mu = \lfloor \frac{k-1}{r} \rfloor = \lceil \frac{k}{r} \rceil - 1$ regardless of whether $r \mid (k-1)$ or not. In addition, if $N = 1$ and $r \nmid (k-1)$, then $\Lambda = 0$. Therefore, the bound directly follows from (10). ∎

Similarly, when $N = \delta - 1$, Theorem 1 is exactly the bound in Lemma 2, first derived in [28].

**Corollary 2:** For any $[n, k, d]_q$ linear code with information $(r, N = \delta - 1, \delta)$-locality,

$$d \leqslant n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil. \tag{16}$$

*Proof:* When $N = \delta - 1$, if $(1 + (\delta-1)(r-1))|(k-1)$, we have

$$\mu(\delta-1) = \frac{(k-1)(\delta-1)}{1 + (\delta-1)(r-1)} = \left\lceil \frac{(k-1)(\delta-1)+1}{(\delta-1)(r-1)+1} \right\rceil - 1,$$

which means

$$\begin{aligned} d &\leqslant n - (k-1) - \mu(\delta-1) \\ &= n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil \end{aligned}$$

according to (10).

When $(1 + (\delta-1)(r-1)) \nmid (k-1)$, it follows from $\Lambda = \lfloor \frac{(k-1)_{1+N(r-1)}-1}{r-1} \rfloor = \lfloor \frac{(k-1)_{1+(\delta-1)(r-1)}-1}{r-1} \rfloor$ that

$$\begin{aligned} \Lambda &\geqslant \frac{k-1}{r-1} - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor \frac{1+(\delta-1)(r-1)}{r-1} - 1 \\ &\geqslant \frac{k-1}{r-1} - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor \frac{1}{r-1} \\ &\quad - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor (\delta-1) - 1 \end{aligned}$$

which is equivalent to

$$\begin{aligned} \Lambda &\geqslant \left\lceil \frac{k-1}{r-1} - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor \frac{1}{r-1} \right\rceil - 1 \\ &\quad - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor (\delta-1) \\ &\geqslant \left\lceil \frac{(k-1)(\delta-1)}{1+(\delta-1)(r-1)} \right\rceil - 1 \\ &\quad - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor (\delta-1). \end{aligned} \tag{17}$$

Therefore, by (10) we have

$$
\begin{aligned}
d \leqslant \;& n - (k-1) - \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor (\delta-1) - \Lambda \\
\leqslant \;& n - (k-1) - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor (\delta-1) \\
& - \left( \left\lceil \frac{(k-1)(\delta-1)}{1+(\delta-1)(r-1)} \right\rceil - 1 \right. \\
& \qquad \left. - \left\lfloor \frac{k-1}{1+(\delta-1)(r-1)} \right\rfloor (\delta-1) \right) \\
= \;& n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{1+(\delta-1)(r-1)} \right\rceil,
\end{aligned}
$$

where the last equality follows from the fact $\left\lceil \frac{(k-1)(\delta-1)}{1+(\delta-1)(r-1)} \right\rceil = \left\lceil \frac{(k-1)(\delta-1)+1}{1+(\delta-1)(r-1)} \right\rceil$ for $(1+(\delta-1)(r-1)) \nmid (k-1)$. This completes the proof. ∎

Generally, we have the following alternative form of Theorem 1.

**Corollary 3:** For any $[n,k,d]_q$ linear code with information $(r,\mathbf{N},\delta)$-locality,

$$
d \leqslant
\begin{cases}
n - k + 1 - \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor (\delta-1), \\
\qquad \text{if } (1+N(r-1)) \mid (k-1), \\[2mm]
n - k + 1 - \left\lceil \frac{\left( \left\lceil \frac{(k-1)N}{1+N(r-1)} \right\rceil - 1 \right)(\delta-1)}{N} \right\rceil, \quad \text{otherwise,}
\end{cases}
$$

where $N = \max(\{N_j : 1 \leqslant j \leqslant k\})$.

*Proof:* This corollary is an immediate result of Theorem 1 by using $\mu = \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor$ and

$$
\Lambda \geqslant \left\lceil \frac{(k-1)N}{1+N(r-1)} \right\rceil - 1 - \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor N
$$

which can be deduced similarly to (17) for $(1+N(r-1)) \nmid (k-1)$. ∎

Corollaries 1–3 tell us that Theorem 1 not only contains the bounds for the cases $N=1$ [9], [20] and $N=\delta-1$ [28], but also provides bounds for other cases. In the following sections, we will prove that these bounds are sometimes tight.

**Remark 5:** Compared with an $[n,k,d]_q$ MDS code, the value $\mu(\delta-1) + \left\lceil \frac{\Lambda(\delta-1)}{N} \right\rceil$ for the case $\langle k-1 \rangle_{1+N(r-1)} \neq 0$ ($\mu(\delta-1)$ for the case $\langle k-1 \rangle_{1+N(r-1)} = 0$, respectively) stands for the least redundancy allowing the code to have information $(r,\mathbf{N},\delta)$-locality according to the Singleton bound, where $\mu = \lfloor \frac{k-1}{1+N(r-1)} \rfloor$ and $\Lambda = \lfloor \frac{\langle k-1 \rangle_{1+N(r-1)} - 1}{r-1} \rfloor$. Thus, for given $r$ and $\delta$, it is easy to check that the smaller $N$ is the larger required redundancy is, when $k-1 \geqslant 1 + N(r-1)$.

## V. LOCALLY REPAIRABLE CODES VIA GABIDULIN CODES

After having proved a bound on the code parameters in the previous section, we turn to providing a construction – Construction A – the first of two. The construction is based on Gabidulin codes with carefully chosen parameters. In particular, the evaluation points for Gabidulin codes need to be chosen, which we first study in Lemma 6. We then give Construction A followed by two main theorems: Theorem 2 finds the locality of the constructed codes, whereas Theorem 3 determines a lower bound on their minimum distance. Several technical lemmas assist in proving the two theorems. Finally, a sequence of corollaries is provided in which specific code parameters are used. In particular, Corollaries 5 and 6 show two families of optimal codes emanating from Construction A.

**Definition 5** ([14]): A polynomial of the form

$$
f(x) = \sum_{i=0}^{k-1} a_i x^{q^i} \tag{18}
$$

with coefficients in an extension field $\mathbb{F}_{q^m}$ of $\mathbb{F}_q$ is called a *q-polynomial* over $\mathbb{F}_{q^m}$. Let $\mathcal{F}(q,m,k)$ denote the set of all possible $q$-polynomials over $\mathbb{F}_{q^m}$ with degree less than $q^k$.

**Lemma 6** ([7]): Let $V = \{v_i : 1 \leqslant i \leqslant n\} \subseteq \mathbb{F}_{q^m}$ and

$$
\mathcal{C} = \{(f(v_1), f(v_2), \cdots, f(v_n)) : f(x) \in \mathcal{F}(q,m,k)\}. \tag{19}
$$

Then,

- $\mathcal{C}$ is an $[n,k]_{q^m}$ linear code if the rank of $V$ over $\mathbb{F}_q$ is greater than or equal to $k$;
- The codeword $C = (f(v_1), f(v_2), \cdots, f(v_n))$ can be recovered by the set of values $\{f(v) : v \in S\}$ if the rank of $S$ over $\mathbb{F}_q$ is greater than or equal to $k$ for any $S \subseteq V$.

In [7], $V$ is required to be linearly independent over $\mathbb{F}_q$ to ensure that $\mathcal{C}$ is an MDS code, which is called a Gabidulin code. In what follows, we intend to propose a construction of codes with information $(r,\mathbf{N},\delta)$-locality by modifying Gabidulin codes. The key difference is to construct a set of vectors $V$, where some elements can be linearly represented by a small number of other elements. Note that a Gabidulin code is based on $f(x)$ in (18), which is a linearized polynomial. In our construction, the linearized property given in (18), and the linear relationship between elements of $V$, will guarantee the desired locality of the code $\mathcal{C}$. More precisely, we have the following construction.

**Construction A:** For any given $\mathbf{N} = (N_1, N_2, \cdots, N_k)$ and $\mathcal{D} = \{d_l^j \geqslant 2 : 1 \leqslant j \leqslant k, 1 \leqslant l \leqslant N_j\}$, let

$$
n = \sum_{1 \leqslant j \leqslant k} \left( 1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2) \right).
$$

We can obtain a linear code by the following steps:

**Step 1**: Select an $[r + d_{\max} - 1, r, d_{\max}]_q$ linear MDS code $\mathcal{C}^*$ whose canonical generator matrix is given as $(I_r, P)$ with $P = (\mathbf{P}_1, \mathbf{P}_2, \cdots, \mathbf{P}_{d_{\max}-1})$, where $d_{\max} = \max(\mathcal{D})$;

**Step 2**: Generate an $(r + d_i^j - 1)$-subset of $\mathbb{F}_{q^m}$, $V_{i,j} = \left\{ v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-2+d_i^j}^{(i,j)} \right\}$ for $1 \leqslant j \leqslant k$ and $1 \leqslant i \leqslant N_j$ satisfying

$$
\begin{aligned}
& \left( v_r^{(i,j)}, v_{r+1}^{(i,j)}, \cdots, v_{r-2+d_i^j}^{(i,j)} \right) \\
& = \left( v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-1}^{(i,j)} \right) (\mathbf{P}_1, \mathbf{P}_2, \cdots, \mathbf{P}_{d_i^j - 1}), \tag{20}
\end{aligned}
$$

where $\left\{ v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-1}^{(i,j)} \right\}$ can be any $r$-subset of $\mathbb{F}_{q^m}$;

**Step 3**: Let $V = \bigcup_{\substack{1 \leqslant j \leqslant k \\ 1 \leqslant i \leqslant N_j}} V_{i,j}$. Construct a code $\mathcal{C}$ with length $|V| \leqslant n$ by means of (19).

Firstly, we have the following theorem for the code generated by Construction A.

**Theorem 2**: For any given positive integers $r$, $k$, $m$ with $r < k$, if $q \geqslant r + d_{\max} - 1$, $V \subseteq \mathbb{F}_{q^m}$, $|V| = n$ and $\text{Rank}(\{v_i : 1 \leqslant i \leqslant k\}) = k$, then the code $\mathcal{C}$ generated by Construction A is an $[n, k]_{q^m}$ linear code with information $(r, \mathbf{N}, \delta)$-locality, where

$$n = \sum_{1 \leqslant j \leqslant k} (1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2)),$$

and

$$\delta = 1 + \min\left(\left\{\sum_{1 \leqslant l \leqslant N_j} (d_l^j - 1) : 1 \leqslant j \leqslant k\right\}\right). \quad (21)$$

*Proof:* It is well known that over $\mathbb{F}_q$ with $q \geqslant r + d_{\max} - 1$, such an MDS code $\mathcal{C}^*$ for Step 1 in Construction A does exist. Since $\text{Rank}(\{v_i : 1 \leqslant i \leqslant k\}) = k$, $|V| = n$, and $V \subseteq \mathbb{F}_{q^m}$, by Lemma 6, we have that the code $\mathcal{C}$ is an $[n, k]_{q^m}$ linear code. This is to say that code symbols $f(v_j)$ for $1 \leqslant j \leqslant k$ can be viewed as the $k$ information symbols.

For $1 \leqslant j \leqslant k$ and $1 \leqslant i \leqslant N_j$, since $(I_r, P)$ is a generator matrix of an $[r + d_{\max} - 1, r, d_{\max}]_q$ MDS code, by (20), we know that any $v \in V_{i,j}$ can be represented as $v = \sum_{v_w^{(i,j)} \in T} e_w^{(i,j,T)} v_w^{(i,j)}$, where $T$ is any $r$-subset of $V_{i,j} \backslash \{v\}$ and $e_w^{(i,j,T)} \in \mathbb{F}_q$. Then, the linearized property of $f(x)$ over $\mathbb{F}_{q^m}$ results in

$$f(v) = f\left(\sum_{v_w^{(i,j)} \in T} e_w^{(i,j,T)} v_w^{(i,j)}\right) = \sum_{v_w^{(i,j)} \in T} e_w^{(i,j,T)} f\left(v_w^{(i,j)}\right)$$

for any $r$-subset $T \subset V_{i,j} \backslash \{v\}$. This is to say the code symbol $f(v)$ can be recovered by $\{f(v_w^{(i,j)}) : v_w^{(i,j)} \in T\}$ for any $r$-subset $T$ of $V_{i,j} \backslash \{v\}$, which means that the code

$$\mathcal{C}|_{V_{i,j}} \triangleq \left\{\left(f(v_j), f\left(v_1^{(i,j)}\right), f\left(v_2^{(i,j)}\right), \cdots, f\left(v_{r+d_i^j-2}^{(i,j)}\right)\right) : \right.$$
$$\left. f(x) \in \mathcal{F}(q, m, k)\right\}$$

is an $[r + d_i^j - 1, r_i^j \geqslant 1, d_i^j]_{q^m}$ linear code for any $1 \leqslant j \leqslant k$ and $1 \leqslant i \leqslant N_j$, where $f(v_j)$ is an information symbol means that $r_i^j \geqslant 1$. Note that

$$|V| = n = \sum_{1 \leqslant j \leqslant k} \left(1 + \sum_{1 \leqslant i \leqslant N_j} \left(r + d_i^j - 2\right)\right)$$

implies that for any $1 \leqslant j \leqslant k$, $V_{i_1,j} \cap V_{i_2,j} = \{v_j\}$ for $1 \leqslant i_1 < i_2 \leqslant k$. Therefore, for any $1 \leqslant j \leqslant k$, $f(v_j)$ has $(r, N_j, \delta)$-locality by Definition 4 and (21), i.e., the code $\mathcal{C}$ has information $(r, \mathbf{N}, \delta)$-locality according to Definition 4. This completes the proof. ∎

Next, we determine the minimum Hamming distance of the code $\mathcal{C}$ generated by Construction A.

**Lemma 7**: For $1 \leqslant j \leqslant k$, denote $V_j = \bigcup_{1 \leqslant i \leqslant N_j} V_{i,j}$. Let

$$V_j \backslash \left\{v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r\right\}$$

be linearly independent over $\mathbb{F}_q$ and have size $1 + N_j(r - 1)$. For any $S \subseteq V_j$, if $\text{Rank}(S) = \tau$, then $|S| \leqslant \tau + \Delta_j(\lfloor \frac{\tau-1}{r-1} \rfloor)$, where $\Delta_j(i) = \sum_{1 \leqslant l \leqslant i} \left(d_l^j - 1\right)$ for any positive integer $1 \leqslant i \leqslant N_j$ and $d_1^j \geqslant d_2^j \geqslant \cdots \geqslant d_{N_j}^j$.

*Proof:* First of all, it is easy to see that the set $V_j$ satisfies the following properties:

P2. Each set $V_{i,j}$ can be spanned by any $r$ of their symbols;

P3. Any two sets of $r$ symbols, respectively from $V_{i_1,j}$ and $V_{i_2,j}$, are linearly dependent;

P4. For any two sets $V_{i_1,j}$ and $V_{i_2,j}$ for $1 \leqslant i_1 < i_2 \leqslant N_j$, $V_{i_1,j} \cap V_{i_2,j} = \{v_j\}$.

This is because P2 is a direct consequence of (20) and the fact that $(I_r, P)$ is a generator matrix of an $[r + d_{\max} - 1, r, d_{\max}]_q$ MDS code. P3 follows immediately from P1, i.e., $V_{i,j}$ for $1 \leqslant i \leqslant N_j$ can be spanned by $v_j$ and any other $r - 1$ symbols. As for P4, clearly $v_j \in V_{i_1,j} \cap V_{i_2,j}$. Assume that there exists an element $v$ such that $\{v_j, v\} \subseteq V_{i_1,j} \cap V_{i_2,j}$. If so, we can find a set $W$ of size $|W| \leqslant 2 + 2(r - 2) = 2r - 2$ containing $v_j$ and $v$ with $|W \cap V_{i_1,j}| = |W \cap V_{i_2,j}| = r$, which by P1, $(V_{i_1,j} \cup V_{i_2,j}) \subseteq \text{Span}(W)$. However, $\text{Rank}(V_{i_1,j} \cup V_{i_2,j}) \geqslant 2r - 1$ since $V_j \backslash \left\{v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r\right\}$ is linearly independent over $\mathbb{F}_q$ and has size $1 + N_j(r - 1)$, a contradiction.

Combining P2-P4 and (23), we know that the hypothesis, namely, $V_j \backslash \left\{v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r\right\}$ is linearly independent over $\mathbb{F}_q$, is equivalent to saying that:

- $v_j$ and any $r - 1$ elements from each set $V_{i,j} \backslash \{v_j\}$; or
- any $r$ elements from one set $V_{i_w,j} \backslash \{v_j\}$ and any $r - 1$ elements from each remaining set $V_{i,j} \backslash \{v_j\}$, with $i \neq i_w$ and $1 \leqslant i \leqslant N_j$, (a total of $1 + N_j(r - 1)$ elements) are linearly independent.

Hence, we have

$\text{Rank}(S)$

$$= \begin{cases} \text{Rank}(\{v_j\}) + \sum_{|S \cap V_{i,j}| \geqslant r} (\text{Rank}(V_{i,j}) - 1) \\ \quad + \sum_{|S \cap V_{i,j}| < r} (|S \cap V_{i,j}| - 1), \quad \text{if } v_j \in S, \\ \sum_{|S \cap V_{i,j}| < r} |S \cap V_{i,j}| + \sum_{\substack{|S \cap V_{i,j}| \geqslant r \\ i \neq i_w}} (\text{Rank}(V_{i,j}) - 1) \\ \quad + \text{Rank}(\{V_{i_w,j}\}), \\ \qquad \text{if } v_j \notin S, \exists i_w, \text{ s.t. } |V_{i_w,j} \cap S| \geqslant r, \\ \sum_{|S \cap V_{i,j}| < r} |S \cap V_{i,j}|, \qquad \text{otherwise,} \end{cases}$$

$$= \begin{cases} 1 + \sum_{|S \cap V_{i,j}| \geqslant r} (r - 1) + \sum_{|S \cap V_{i,j}| < r} (|S \cap V_{i,j}| - 1), \\ \qquad \text{if } v_j \in S, \\ 1 + \sum_{|S \cap V_{i,j}| \geqslant r} (r - 1) + \sum_{|S \cap V_{i,j}| < r} |S \cap V_{i,j}|, \qquad (22) \\ \qquad \text{if } v_j \notin S \text{ and } \exists i_w, \text{ s.t. } |V_{i_w,j} \cap S| \geqslant r, \\ |S|, \qquad \text{otherwise,} \end{cases}$$

where for $v_j \notin S$, we set $V_{i_w,j} = \emptyset$ if $|S \cap V_{i_w,j}| < r$ for all $1 \leqslant i_w \leqslant N_j$, otherwise we choose a set $V_{i_w,j}$ with $|S \cap V_{i_w,j}| \geqslant r$.

It follows from P2 and the fact that $|V_{i,j}| = r + d_i^j - 1$ for all $1 \leqslant i \leqslant N_j$, that

$$|S| \leqslant \begin{cases} 1 + \sum_{1 \leqslant i \leqslant N_j} |(S \cap V_{i,j}) \setminus \{v_j\}|, & \text{if } v_j \in S, \\ \sum_{1 \leqslant i \leqslant N_j} |(S \cap V_{i,j}) \setminus \{v_j\}|, & \text{otherwise,} \end{cases}$$

$$\leqslant \begin{cases} 1 + \sum_{|S \cap V_{i,j}| \geqslant r} (r + d_i^j - 2) \\ \quad + \sum_{|S \cap V_{i,j}| < r} (|S \cap V_{i,j}| - 1), & \text{if } v_j \in S, \\ \sum_{|S \cap V_{i,j}| \geqslant r} (r + d_i^j - 2) \\ \quad + \sum_{|S \cap V_{i,j}| < r} |S \cap V_{i,j}|, & \text{otherwise.} \end{cases} \tag{23}$$

Finally, comparing (22) with (23), we have

$$|S| \leqslant \text{Rank}(S) + \sum_{|S \cap V_{i,j}| \geqslant r} \left( d_i^j - 1 \right)$$
$$\leqslant \text{Rank}(S) + \Delta_j(M),$$

where $M = |\{V_{i,j} : |S \cap V_{i,j}| \geqslant r, 1 \leqslant i \leqslant N_j\}|$ and in the second inequality we use $\sum_{|S \cap V_{i,j}| \geqslant r} \left( d_i^j - 1 \right) \leqslant \Delta_j(M)$ from the assumption $d_1^{(\mathbf{I}_j)} \geqslant d_2^{(\mathbf{I}_j)} \geqslant \cdots \geqslant d_{N_j}^{(\mathbf{I}_j)}$. Obviously, $M \leqslant \lfloor \frac{\tau-1}{r-1} \rfloor$ by (22), which completes the proof. ∎

Consider $\left( d_1^1, d_2^1, \cdots, d_{N_1}^1, d_1^2, \cdots, d_{N_k}^k \right)$ and reorder its elements as $(d_1, d_2, \cdots, d_u)$ such that $d_1 \geqslant d_2 \geqslant \cdots \geqslant d_u$ and $u = \sum_{1 \leqslant j \leqslant k} N_j$. Define $\Delta(t) = \sum_{1 \leqslant j \leqslant t} (d_j - 1)$ for $1 \leqslant t \leqslant u$.

**Lemma 8:** For $1 \leqslant j \leqslant k$, let $V_j = \bigcup_{1 \leqslant i \leqslant N_j} V_{i,j}$, $V = \bigcup_{1 \leqslant j \leqslant k} V_j$, and $N = \max(\{N_j : 1 \leqslant j \leqslant k\})$.

1. If
$$V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$$
is linearly independent over $\mathbb{F}_q$ and has size $1 + N_j(r-1)$ for $1 \leqslant j \leqslant k$, then $\sum_{1 \leqslant j \leqslant k} \text{Rank}(V' \cap V_j) \geqslant k$ for any $(k + \Delta(\lfloor \frac{(k-1)N}{(r-1)N+1} \rfloor))$-subset $V'$ of $V$;

2. Furthermore, if
$$V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$$
is linearly independent over $\mathbb{F}_q$ and has size $\sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$, then $\text{Rank}(V') = \text{Rank}(\sum_{1 \leqslant i \leqslant k} V' \cap V_j) \geqslant k$ for any $\left( k + \Delta \left( \lfloor \frac{(k-1)N}{(r-1)N+1} \rfloor \right) \right)$-subset $V'$ of $V$.

*Proof:* For the first part, we assume to the contrary that there exists a set $V'$ with

$$|V'| = k + \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right) \tag{24}$$

but

$$\sum_{1 \leqslant i \leqslant k} \text{Rank}(S_j) \leqslant k - 1, \tag{25}$$

where we set $S_j = V' \cap V_j$ for $1 \leqslant j \leqslant k$.

The fact that
$$\left| V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \right| = 1 + N_j(r-1)$$

for each $1 \leqslant j \leqslant k$ means

$$\left\lfloor \frac{\text{Rank}(S_j) - 1}{r - 1} \right\rfloor = \left\lfloor \frac{\text{Rank}(S_j)N}{(r-1)N+1} \right\rfloor,$$

since $\text{Rank}(S_j) \leqslant \text{Rank}(V_j) = 1 + N_j(r-1) \leqslant 1 + N(r-1)$ and $\lfloor \frac{t-1}{r-1} \rfloor = \lfloor \frac{tN}{(r-1)N+1} \rfloor$ for any positive integer $t \leqslant 1 + N(r-1)$. Thus, by Lemma 7 and (25),

$$|V'| \leqslant \sum_{1 \leqslant j \leqslant k} |S_j|$$
$$\leqslant \sum_{1 \leqslant j \leqslant k} \left( \text{Rank}(S_j) + \Delta_j \left( \left\lfloor \frac{\text{Rank}(S_j) - 1}{r - 1} \right\rfloor \right) \right)$$
$$= \sum_{1 \leqslant j \leqslant k} \left( \text{Rank}(S_j) + \Delta_j \left( \left\lfloor \frac{\text{Rank}(S_j)N}{(r-1)N+1} \right\rfloor \right) \right)$$
$$\leqslant k - 1 + \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right),$$

where the last inequality follows from $d_1 \geqslant d_2 \geqslant \cdots \geqslant d_u$, i.e., for $1 \leqslant a_j \leqslant N_j$ and $1 \leqslant j \leqslant k$,

$$\sum_{1 \leqslant j \leqslant k} \Delta_j(a_j) \leqslant \max_{\substack{\Gamma \subseteq [u] \\ |\Gamma| = \sum_{1 \leqslant j \leqslant k} a_j}} \left( \sum_{\tau \in \Gamma} (d_\tau - 1) \right)$$
$$= \Delta \left( \sum_{1 \leqslant j \leqslant k} a_j \right),$$

which contradicts (24). Thus, the desired result follows.

For the second part, the fact that

$$\left| V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \right|$$
$$= \sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$$

means that

$$\left| V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \right| = 1 + N_j(r-1)$$

for any $1 \leqslant j \leqslant k$. Since $V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$, we have

$$\text{Span}(V) = \bigoplus_{1 \leqslant j \leqslant k} \text{Span}(V_j), \tag{26}$$

and $V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is also linearly independent over $\mathbb{F}_q$, where "$\bigoplus$" denotes the direct sum of linear spaces. According to (26) and the result of the first part,

$$\text{Rank}(V') = \text{Rank} \left( \sum_{1 \leqslant j \leqslant k} V' \cap V_j \right)$$
$$= \sum_{1 \leqslant j \leqslant k} \text{Rank}(V' \cap V_j) \geqslant k,$$

for any $V' \subseteq V$ with $|V'| = k + \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$. ∎

Based on Lemmas 7 and 8, we are able to prove our result on the minimum Hamming distance.

**Theorem 3:** For $1 \leqslant j \leqslant k$, let $V_j = \bigcup_{1 \leqslant i \leqslant N_j} V_{i,j}$ and $V = \bigcup_{1 \leqslant j \leqslant k} V_j$. If $q \geqslant r + d_{\max} - 1$ and

$$V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \subseteq \mathbb{F}_{q^m}$$

is linearly independent over $\mathbb{F}_q$ and has size $\sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$, then the code $\mathcal{C}$ generated by Construction A is an $[n, k, d]_{q^m}$ linear code $\mathcal{C}$ with information $(r, \mathbf{N}, \delta)$-locality, where $\mathbf{N} = (N_1, N_2, \cdots, N_k)$, $n = \sum_{1 \leqslant j \leqslant k} (1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2))$ and $d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$.

*Proof:* According to Theorem 2, it suffices to show that $|V| = n = \sum_{1 \leqslant j \leqslant k} (1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2))$ and $d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$. By P4 in the proof of Lemma 7, the facts that $V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \subseteq \mathbb{F}_{q^m}$ is linearly independent over $\mathbb{F}_q$ and

$$\left| V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \right|$$
$$= \sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$$

mean that $V_{i_1,j} \cap V_{i_2,j} = \{v_j\}$ for $1 \leqslant i_1 < i_2 \leqslant N_j, 1 \leqslant j \leqslant k$ and $V_{j_1} \cap V_{j_2} = \emptyset$ for $1 \leqslant j_1 < j_2 \leqslant k$, i.e.,

$$|V| = \sum_{1 \leqslant j \leqslant k} |V_j| = \sum_{1 \leqslant j \leqslant k} \left( 1 + \sum_{1 \leqslant i \leqslant N_j} \left( r + d_i^j - 2 \right) \right) = n.$$

For the minimum Hamming distance $d$ of $\mathcal{C}$, we have

$$d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$$

according to Lemma 6 and Lemma 8-2. This completes the proof. ∎

**Corollary 4:** Let $\mathbf{N} = (N_1, N_2, \cdots, N_k)$ be a sequence of positive integers, $\mathcal{D} = \{d_i^j : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j\}$ and $N = \max(\{N_j : 1 \leqslant j \leqslant k\})$. Denote

$$\delta = 1 + \min \left( \left\{ \sum_{1 \leqslant l \leqslant N_j} \left( d_l^j - 1 \right) : 1 \leqslant j \leqslant k \right\} \right), \quad (27)$$

and $d_{\max} = \max(\mathcal{D})$. For any given positive integers $r, k, m$ with $r < k$, if $m \geqslant k((r-1)N+1)$, $q \geqslant r + d_{\max} - 1$, then Construction A can generate an $[n, k, d]_{q^m}$ linear code $\mathcal{C}$ with information $(r, \mathbf{N}, \delta)$-locality, where $n = \sum_{1 \leqslant j \leqslant k} (1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2))$ and $d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$.

*Proof:* Since Construction A has no restriction on the set $V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$, the hypothesis $m \geqslant k((r-1)N+1) \geqslant \sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$ implies that we can select the set

$$V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\} \subseteq \mathbb{F}_{q^m}$$

to be linearly independent over $\mathbb{F}_q$ with size $\sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$ in Step 2, Construction A. For instance, we can let it be a $\sum_{1 \leqslant j \leqslant k} (1 + N_j(r-1))$-subset of a base for $\mathbb{F}_{q^m}$ over $\mathbb{F}_q$. Now the corollary follows directly from Theorem 3. ∎

In particular, we have the following two specific optimal constructions.

**Corollary 5:** Let $d_1 = d_2 = \cdots = d_u = 2$, $N_j = \delta - 1$ for $1 \leqslant j \leqslant k$, $m \geqslant k((r-1)(\delta-1)+1)$, and $q \geqslant 2$. If $V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$, then the code $\mathcal{C}$ generated by Construction A is an optimal $[n, k, d]_{q^m}$ linear code with information $(r, \delta-1, \delta)$-locality with respect to the bound in Lemma 2, where $n = k(1 + r(\delta-1))$ and $d = n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$.

*Proof:* For the case $d_1 = d_2 \cdots = d_u = 2 = d_{\max}$, to make sure that $[r+1, r, 2]_q$ linear code $\mathcal{C}^*$ exists for Step 1 of Construction A we only need $q \geqslant 2$ rather than $q \geqslant r + d_{\max} - 1$. Thus, by Theorem 3, the code $\mathcal{C}$ is an $[n, k, d]_{q^m}$ linear code with information $(r, \delta-1, \delta)$-locality and

$$d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$$
$$= n - k + 1 - \left\lfloor \frac{(k-1)(\delta-1)}{(r-1)(\delta-1)+1} \right\rfloor$$
$$= n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil,$$

where $n = k(1 + r(\delta-1))$. Recall that by Lemma 2, $d \leqslant n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$. Then $d = n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$ and the code $\mathcal{C}$ is an optimal linear code with information $(r, \delta-1, \delta)$-locality with respect to the bound in Lemma 2. ∎

**Corollary 6:** Let $d_1 = d_2 = \cdots = d_u = d^* > 2$, $N_j = N$ for $1 \leqslant j \leqslant k$, $\delta - 1 = N(d^* - 1)$, $q \geqslant r + d^* - 1$ and $m \geqslant k((r-1)N+1)$. If $V \setminus \left\{ v_t^{(i,j)} : 1 \leqslant j \leqslant k, \ 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$ and has size $k((r-1)N+1)$, then the code $\mathcal{C}$ generated by Construction A is an optimal $[n, k, d]_{q^m}$ linear code with information $(r, N, \delta)$-locality with respect to the bound in Corollary 3, where $n = k(1 + N(r + d^* - 2))$ and $d = n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1)$.

*Proof:* According to Theorem 3, the code $\mathcal{C}$ is an $[n, k, d]_{q^m}$ linear code with information $(r, N, N(d^* - 1)+1)$-locality, where $n = k(1 + N(r + d^* - 2))$ and

$$d \geqslant n - k + 1 - \Delta \left( \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor \right)$$
$$= n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1). \quad (28)$$

In the case $N_j = N$ for $1 \leqslant j \leqslant k$ and $\delta - 1 = N(d^* - 1)$, by Corollary 3,

$$d \leqslant \begin{cases} n - k + 1 - \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor (\delta - 1), \\ \qquad \text{if } (1 + N(r-1)) \mid (k-1), \\ n - k + 1 - \left\lceil \frac{\left( \left\lceil \frac{(k-1)N}{1+N(r-1)} \right\rceil - 1 \right)(\delta - 1)}{N} \right\rceil, \text{ otherwise,} \end{cases}$$
$$= n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1).$$

Therefore, $d = n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1)$ and the code $\mathcal{C}$ is an optimal linear code with $(r, N, \delta)$-locality with respect to the bound in Corollary 3. ∎

We conclude this section with four remarks for locally repairable codes by Construction A.

**Remark 6:** In [28], Wang and Zhang showed the existence of optimal $[n, k, d]_q$ linear codes with information $(r, \delta-1, \delta)$-locality via the Sparse Zero Lemma [6], when $n \geqslant k(r(\delta-1)+1)$ and $q > 1 + \binom{n}{k+\sigma}$ with $\sigma = \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$. However, to the best of our knowledge, no explicit construction has achieved the bound in Lemma 2. Thus, Construction A seems to be the first explicit construction that can yield optimal locally repairable codes with respect to the bound in Lemma 2.

**Remark 7:** If $N = 1$, then Construction A is exactly the one introduced in [21] for optimal locally repairable codes with respect to the bound in Lemma 1. Thus, Construction A can be viewed as a generalization of the one in [21] for the codes with multiple disjoint repair sets.

**Remark 8:** Construction A and Corollaries 5 and 6 also show that the bound in Theorem 1 is tight for some cases.

**Remark 9:** The fact that $\bigcup_{\substack{1 \leqslant i \leqslant k \\ 1 \leqslant j \leqslant N_i}} V_{j,i} = [n]$ where the sets $V_{j,i} \setminus \{v_i\}$ for $1 \leqslant i \leqslant k$ and $1 \leqslant j \leqslant N_i$ correspond to the repair sets for the $k$ information symbols implies that all the $n$ code symbols have locality $r$. However, besides the $k$ information symbols corresponding to $\{v_1, v_2, \cdots, v_k\}$, it is not clear that the other code symbols also have multiple repair sets or their repair sets would tolerate overall $\delta - 1$ erasures. In fact, for all symbol locality, generally how to construct an optimal locally repairable code with multiple repair sets is still an open problem. For further discussion on this problem the reader is referred to [27].

## VI. LOCALLY REPAIRABLE CODES VIA LINEARIZED REED-SOLOMON CODES

Inspired by the constructions in [17] for maximal recoverable codes (or Partial MDS codes), we also employ linearized Reed-Solomon codes to reduce the size of the finite field required for optimal locally repairable codes. This section briefly describes linearized Reed-Solomon codes in Definition 7, citing [17] in Lemma 9. We then give Construction B which replaces Gabidulin codes with linearized Reed-Solomon codes as the building block. Then Theorem 4 provides an analysis of the locality and minimum distance of the constructed code. As in the previous section, two corollaries present specific parameter choices for the construction: Corollaries 7 and 8 give two families of optimal codes emanating from Construction B.

We start by recalling some necessary definitions for linearized Reed-Solomon codes. For positive integers $M$ and $g$, let $L = (L_1, L_2, \ldots, L_g)$, $M = L_1 + L_2 + \cdots + L_g$ and $1 \leqslant L_i \leqslant m$. Let $q$ be a prime power with $q - 1 \geqslant g$. Define $\sigma : \mathbb{F}_{q^m} \to \mathbb{F}_{q^m}$ as $\sigma(\alpha) \triangleq \alpha^q$. We first recall the definition of a linear operator over a finite field as in [16].

**Definition 6:** For any $\alpha \in \mathbb{F}_{q^m}$ and $i \in \mathbb{N}$, define $\mathrm{Norm}_i(\alpha) \triangleq \sigma^{i-1}(\alpha) \cdots \sigma(\alpha)\alpha$. The $\mathbb{F}_q$-linear operator $\Psi^i_\alpha : \mathbb{F}_{q^m} \to \mathbb{F}_{q^m}$ is defined by

$$\Psi^i_\alpha(\beta) = \sigma^i(\beta) \, \mathrm{Norm}_i(\alpha). \tag{29}$$

**Definition 7:** Let $\gamma$ be a primitive element of $\mathbb{F}_{q^m}$ and let $\mathcal{B} = \{\beta_1, \beta_2, \cdots, \beta_m\}$ be a basis of $\mathbb{F}_{q^m}$ over $\mathbb{F}_q$. For $1 \leqslant i \leqslant g$ and $k \in \mathbb{N}$, define the matrices

$$D_i^{(k)} = \begin{pmatrix} \beta_1 & \beta_2 & \cdots & \beta_{L_i} \\ \Psi^1_{\gamma^{i-1}}(\beta_1) & \Psi^1_{\gamma^{i-1}}(\beta_2) & \cdots & \Psi^1_{\gamma^{i-1}}(\beta_{L_i}) \\ \vdots & \vdots & \ddots & \vdots \\ \Psi^{k-1}_{\gamma^{i-1}}(\beta_1) & \Psi^{k-1}_{\gamma^{i-1}}(\beta_2) & \cdots & \Psi^{k-1}_{\gamma^{i-1}}(\beta_{L_i}) \end{pmatrix}.$$

The *linearized Reed-Solomon code* with dimension $k$, primitive element $\gamma$, and basis $\mathcal{B}$ is the linear code $\mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma) \subseteq \mathbb{F}_{q^m}^n$ with generator matrix

$$D = \left( D_1^{(k)}, D_2^{(k)}, \cdots, D_g^{(k)} \right)_{k \times M}. \tag{30}$$

Let $\mathrm{Diag}(W_1, W_2, \cdots, W_g)$ denote the block-diagonal matrix, whose main-diagonal blocks are $W_1, W_2, \cdots, W_g$, i.e.,

$$\mathrm{Diag}(W_1, W_2, \cdots, W_g) = \begin{pmatrix} W_1 & 0 & \cdots & 0 \\ 0 & W_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & W_g \end{pmatrix}.$$

The following property is introduced in [17].

**Lemma 9** ([17])**:** Let $\mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma)$ be the $[M, k]_{q^m}$ linearized Reed-Solomon code in Definition 7 with $M = L_1 + L_2 + \cdots + L_g$. Then for all integers $n_i \geqslant 1$ and all matrices $W_i \in \mathbb{F}_q^{L_i \times n_i}$, for $1 \leqslant i \leqslant g$, satisfying

$$\sum_{1 \leqslant i \leqslant g} \mathrm{Rank}(W_i) \geqslant k, \tag{31}$$

there exists a decoder

$$\mathrm{Dec} : \mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma) \, \mathrm{Diag}(W_1, W_2, \cdots, W_g) \to \mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma)$$

such that

$$\mathrm{Dec}(C \, \mathrm{Diag}(W_1, W_2, \cdots, W_g)) = C \quad \text{for any } C \in \mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma),$$

where

$$\begin{aligned} &\mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma) \, \mathrm{Diag}(W_1, W_2, \cdots, W_g) \\ &\triangleq \{C \, \mathrm{Diag}(W_1, W_2, \cdots, W_g) : C \in \mathcal{C}^\sigma_{L,k}(\mathcal{B}, \gamma)\}. \end{aligned}$$

By replacing the Gabidulin code with a linearized Reed-Solomon code in Construction A, we get the following construction.

**Construction B:** For any given $\mathbf{N} = (N_1, N_2, \cdots, N_k)$ and $\mathcal{D} = \{d_l^j \geqslant 2 : 1 \leqslant j \leqslant k, 1 \leqslant l \leqslant N_j\}$, let

$$n_j = 1 + \sum_{1 \leqslant i \leqslant N_j} (r + d_i^j - 2),$$

for $1 \leqslant j \leqslant k$, and define $n = \sum_{1 \leqslant j \leqslant k} n_j$. Let $g = k$, $L_j = 1 + N_j(r - 1)$, $M = \sum_{1 \leqslant j \leqslant k} L_i$, for $1 \leqslant j \leqslant k$. Assume $m \geqslant L_j$ for $1 \leqslant j \leqslant k$. We can obtain a linear code by the following steps:

**Step 1:** Select an $[r + d_{\max} - 1, r, d_{\max}]_q$ linear MDS code $\mathcal{C}^*$ whose canonical generator matrix is given as $(I_r, P)$ with $P = (\mathbf{P}_1, \mathbf{P}_2, \cdots, \mathbf{P}_{d_{\max}-1})$, where $d_{\max} = \max(\mathcal{D})$;

**Step 2**: Generate an $(r + d_i^j - 1)$-subset of $\mathbb{F}_{q^m}$, $V_{i,j} = \left\{ v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-2+d_i^j}^{(i,j)} \right\}$ for $1 \leqslant j \leqslant k$ and $1 \leqslant i \leqslant N_j$ satisfying

$$\left( v_r^{(i,j)}, v_{r+1}^{(i,j)}, \cdots, v_{r-2+d_i^j}^{(i,j)} \right)$$
$$= \left( v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-1}^{(i,j)} \right) \left( \mathbf{P}_1, \mathbf{P}_2, \cdots, \mathbf{P}_{d_i^j - 1} \right), \quad (32)$$

where $\left\{ v_j, v_1^{(i,j)}, v_2^{(i,j)}, \cdots, v_{r-1}^{(i,j)} \right\}$ can be any $r$-subset of $\mathbb{F}_{q^m}$. Then, based on (32), for each $1 \leqslant j \leqslant k$, an $L_j \times n_j$ matrix $A_j$, can be uniquely determined as follows

$$\mathbf{V}_j = \left( v_j, v_1^{(1,j)}, v_2^{(1,j)}, \cdots, v_{r+d_1^j-2}^{(1,j)}, v_1^{(2,j)}, \right.$$
$$\left. \cdots, v_{r+d_2^j-2}^{(2,j)}, \cdots, v_{r+d_{N_j}^j-2}^{(N_j,j)} \right)$$
$$= \left( v_j, v_1^{(1,j)}, v_2^{(1,j)}, \cdots, v_{r-1}^{(1,j)}, v_1^{(2,j)}, \right.$$
$$\left. \cdots, v_{r-1}^{(2,j)}, \cdots, v_{r-1}^{(N_j,j)} \right) A_j. \quad (33)$$

**Step 3**: Let $D$ be the generator matrix of the $[M, k]_{q^m}$ linearized Reed-Solomon code $\mathcal{C}_{L,k}^\sigma(\mathcal{B}, \gamma)$. Construct a code $\mathcal{C}$ with length $n$ over $\mathbb{F}_{q^m}$ by the generator matrix $G = D \operatorname{Diag}(A_1, A_2, \cdots, A_k)$, i.e.,

$$\mathcal{C} = \mathcal{C}_{L,k}^\sigma(\mathcal{B}, \gamma) \operatorname{Diag}(A_1, A_2, \cdots, A_k)$$
$$\triangleq \{ C \operatorname{Diag}(A_1, A_2, \cdots, A_k) : C \in \mathcal{C}_{L,k}^\sigma(\mathcal{B}, \gamma) \}.$$

Note from (33) that

P5. For $1 \leqslant j \leqslant k$, if $V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$, then

$$\operatorname{Rank}(\mathbf{V}_j(S)) = \operatorname{Rank}(A_j(S)), \quad (34)$$

where for $S = \{s_1, s_2, \cdots, s_t\} \subseteq [L_j]$,

$$\mathbf{V}_j = (v_{j,1}, v_{j,2}, \cdots, v_{j,L_j}) \in \mathbb{F}_{q^m}^{L_j},$$

and

$$A_j = (A_{1,1}, A_{1,1}, \cdots, A_{1,L_j}),$$

we define

$$\mathbf{V}_j(S) \triangleq (v_{j,s_1}, v_{j,s_2}, \cdots, v_{j,s_t}),$$

and

$$A_j(S) \triangleq (A_{j,s_1}, A_{j,s_1}, \cdots, A_{j,s_t}).$$

Then, applying it to Lemma 9, the requirement on the rank of submatrix $A_j(S)$ can be transformed to the rank of the corresponding subset $\mathbf{V}_j(S)$. Immediately, using Lemma 8, we get the following result.

**Theorem 4**: For $1 \leqslant j \leqslant k$, let $V_j = \bigcup_{1 \leqslant i \leqslant N_j} V_{i,j}$, $V = \bigcup_{1 \leqslant j \leqslant k} V_j$, and $m = \max_{1 \leqslant j \leqslant k}(1 + N_j(r-1))$, where

$$V_j = \left\{ v_j, v_1^{(1,j)}, v_2^{(1,j)}, \cdots, v_{r+d_1^j-2}^{(1,j)}, v_1^{(2,j)}, \right.$$
$$\left. \cdots, v_{r+d_2^j-2}^{(2,j)}, \cdots, v_{r+d_{N_j}^j-2}^{(N_j,j)} \right\}.$$

If $q \geqslant \max\{k + 1, r + d_{max} - 1\}$ and

$$V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$$

is linearly independent over $\mathbb{F}_q$ and has size $1 + N_j(r-1)$ for $1 \leqslant j \leqslant k$, then the code $\mathcal{C}$ generated by Construction B is an $[n, k, d]_{q^m}$ linear code $\mathcal{C}$ with information $(r, \mathbf{N}, \delta)$-locality, where $n = \sum_{1 \leqslant j \leqslant k}(1 + \sum_{1 \leqslant i \leqslant N_j}(r + d_i^j - 2))$ and $d \geqslant n - k + 1 - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right)$ and $\delta = 1 + \min\left(\left\{ \sum_{1 \leqslant l \leqslant N_j}(d_l^j - 1) : 1 \leqslant j \leqslant k \right\}\right)$.

*Proof:* The fact that $q \geqslant r + d_{\max} - 1$ guarantees the existence of the MDS code $\mathcal{C}^*$ over $\mathbb{F}_q$ for Step 1 in Construction B. Further, the facts $m = \max_{1 \leqslant j \leqslant k}(1 + N_j(r-1))$ and $q \geqslant k + 1$ imply that the linearized Reed-Solomon code for Step 3 in Construction B exists. By Construction B, $\mathcal{C}$ is an $[n, k]_{q^m}$ code.

For the convenience of discussion, we index the codeword $C \in \mathcal{C}_{L,k}^\sigma(\mathcal{B}, \gamma)$ as

$$C = \left( c_1, c_1^{(1,1)}, \cdots, c_{r-1}^{(1,1)}, c_1^{(2,1)}, \cdots, c_{r-1}^{(2,1)}, \cdots, c_1^{(N_1,1)}, \cdots, \right.$$
$$\left. c_{r-1}^{(N_1,1)}, c_2, c_1^{(1,2)}, \cdots, c_{r-1}^{(N_2,2)}, \cdots c_k, c_1^{(1,k)}, \cdots c_{r-1}^{(N_k,k)} \right)$$

and

$$C' = \left( c_1, c_1^{(1,1)}, \cdots, c_{r+d_1^1-2}^{(1,1)}, c_1^{(2,1)}, \cdots, c_{r+d_2^1-2}^{(2,1)}, \cdots, c_1^{(N_1,1)} \right.$$
$$\cdots, c_{r+d_{N_1}^1-2}^{(N_1,1)}, c_2, c_1^{(1,2)}, \cdots, c_{r+d_{N_2}^2-2}^{(N_2,2)}, \cdots c_k, c_1^{(1,k)},$$
$$\left. \cdots c_{r+d_{N_k}^k-2}^{(N_k,k)} \right)$$

for $C' = C \operatorname{Diag}(A_1, A_2, \cdots, A_k) \in \mathcal{C}$.

Firstly we claim that $c_j$ for $1 \leqslant j \leqslant k$ can be regarded as information symbols. By (34), we have $k = \sum_{1 \leqslant j \leqslant k} \operatorname{Rank}((v_j)) = \sum_{1 \leqslant j \leqslant k} \operatorname{Rank}(A_{j,1})$. According to Lemma 9, this means that the code symbols $c_j$ ($1 \leqslant j \leqslant k$) are able to recover the whole codeword $C$ and then $C'$. Thus, the claim follows.

Next, we prove the locality of the code symbol $c_j$ for $1 \leqslant j \leqslant k$. For each $1 \leqslant j \leqslant k$ and $1 \leqslant t \leqslant N_j$, equations (32) and (33) mean that

$$\left( c_j, c_1^{(i,j)}, c_2^{(i,j)}, \cdots, c_{r+d_i^j-2}^{(i,j)} \right)$$
$$= \left( c_j, c_1^{(i,j)}, c_2^{(i,j)}, \cdots, c_{r-1}^{(i,j)} \right) \left( I_r, \mathbf{P}_1, \mathbf{P}_2, \ldots, \mathbf{P}_{d_i^j-1} \right).$$

Hence, the punctured code

$$\mathcal{C}_{V_{i,j}} \triangleq \left\{ \left( c_j, c_1^{(i,j)}, c_2^{(i,j)}, \cdots, c_{r+d_i^j-2}^{(i,j)} \right) : C' \in \mathcal{C} \right\}$$

is an $[r + d_i^j - 1, r_i^j \geqslant 1, d_i^j]_{q^m}$ linear code, where $r_i^j \geqslant 1$ follows by the fact that $c_j$ for $1 \leqslant j \leqslant k$ is an information symbol. Therefore, by Definition 4, the code symbol $c_j$ for $1 \leqslant j \leqslant k$ has $(r, N_j, \delta)$-locality, i.e., the code $\mathcal{C}$ generated by Construction B has information $(r, \mathbf{N}, \delta)$-locality, where $\delta = 1 + \min\left(\left\{ \sum_{1 \leqslant l \leqslant N_j}(d_l^j - 1) : 1 \leqslant j \leqslant k \right\}\right)$.

As for the minimum Hamming distance $d$ of $\mathcal{C}$, assume that erasure pattern is $E$ with $|E| \leqslant n - k - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right)$

and $S_j \subseteq [L_j]$ is the set of the indices for the elements of $V_j \setminus E$ over $\mathbf{V}_j$ for $1 \leqslant j \leqslant k$. Recall that $V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$ and has size $1 + N_j(r - 1)$ for $1 \leqslant j \leqslant k$. According to Lemma 8-1, we have $\sum_{1 \leqslant j \leqslant k} \mathrm{Rank}(V_j \setminus E) \geqslant k$. Immediately, it follows from (34) that $\sum_{1 \leqslant j \leqslant k} \mathrm{Rank}(A_j(S_j)) = \sum_{1 \leqslant j \leqslant k} \mathrm{Rank}(\mathbf{V}_j(S_j)) = \sum_{1 \leqslant j \leqslant k} \mathrm{Rank}(V_j \setminus E) \geqslant k$. That is, any erasure pattern $E$ with $|E| \leqslant n - k - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right)$ can be recovered by Lemma 9. Therefore,

$$d \geqslant n - k + 1 - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right),$$

which completes the proof. ■

In what follows, we discuss two specific settings in which Construction B yields optimal codes.

**Corollary 7:** Let $d_1 = d_2 = \cdots = d_u = 2$, $N_j = \delta - 1$ for $1 \leqslant j \leqslant k$, $m \geqslant (r-1)(\delta-1)+1$ and $q \geqslant k+1$. If $V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$ is linearly independent over $\mathbb{F}_q$ and has size $(r-1)(\delta-1)+1$, then the code $\mathcal{C}$ generated by Construction B is an optimal $[n, k, d]_{q^m}$ linear code with information $(r, \delta-1, \delta)$-locality with respect to the bound in Lemma 2, where $n = k(1 + r(\delta - 1))$ and $d = n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$.

*Proof:* By Theorem 4, the code $\mathcal{C}$ is an $[n, k, d]_{q^m}$ linear code with information $(r, \delta-1, \delta)$-locality and

$$d \geqslant n - k + 1 - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right)$$
$$= n - k + 1 - \left\lfloor \frac{(k-1)(\delta-1)}{(r-1)(\delta-1)+1} \right\rfloor$$
$$= n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil,$$

where $n = k(1 + r(\delta - 1))$. Recall that by Lemma 2, $d \leqslant n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$. Then $d = n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$ and the code $\mathcal{C}$ is an optimal linear code with information $(r, \delta-1, \delta)$-locality with respect to the bound in Lemma 2. ■

**Corollary 8:** Let $d_1 = d_2 = \cdots = d_u = d^* > 2$, $N_j = N$ for $1 \leqslant j \leqslant k$, $\delta - 1 = N(d^* - 1)$, $q \geqslant \max\{r + d^* - 1, k + 1\}$ and $m \geqslant (r-1)N + 1$. If

$$V_j \setminus \left\{ v_t^{(i,j)} : 1 \leqslant i \leqslant N_j, t \geqslant r \right\}$$

is linearly independent over $\mathbb{F}_q$ and has size $(r-1)N + 1$, then the code $\mathcal{C}$ generated by Construction B is an optimal $[n, k, d]_{q^m}$ linear code with information $(r, N, \delta)$-locality with respect to the bound in Corollary 3, where $n = k(1 + N(r + d^* - 2))$ and $d = n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1)$.

*Proof:* According to Theorem 4, the code $\mathcal{C}$ is an $[n, k, d]_{q^m}$ linear code with information $(r, N, N(d^*-1)+1)$-locality, where $n = k(1 + N(r + d^* - 2))$ and

$$d \geqslant n - k + 1 - \Delta\left(\left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor\right)$$
$$= n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1).$$

In the case $N_j = N$ for $1 \leqslant j \leqslant k$ and $\delta - 1 = N(d^* - 1)$, by Corollary 3,

$$d \leqslant \begin{cases} n - k + 1 - \left\lfloor \frac{k-1}{1+N(r-1)} \right\rfloor (\delta - 1), \\ \qquad\qquad \text{if } (1 + N(r-1)) \mid (k-1), \\ n - k + 1 - \left\lceil \frac{\left(\left\lceil \frac{(k-1)N}{1+N(r-1)} \right\rceil - 1\right)(\delta-1)}{N} \right\rceil, \\ \qquad\qquad \text{otherwise,} \end{cases}$$

$$= n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1).$$

Therefore, $d = n - k + 1 - \left\lfloor \frac{(k-1)N}{(r-1)N+1} \right\rfloor (d^* - 1)$ and the code $\mathcal{C}$ is an optimal linear code with information $(r, N, \delta)$-locality with respect to the bound in Corollary 3. ■

We conclude this section by an illustrative example for an optimal locally repairable code generated by Construction B.

**Example 1:** Let $k = 3$, $r = 2$, $\delta = 3$, and $N = 2$. Set $n = 15$, $L_1 = L_2 = L_3 = 3$, and $M = L_1 + L_2 + L_3 = 9$. Note that in this case $d_i^j = 2$ for $1 \leqslant j \leqslant 3$ and $1 \leqslant i \leqslant 2$. Thus, the required field size for the linearized Reed-Solomon code is $q \geqslant 4$ and $m \geqslant 3$. Apply the primitive polynomial $f(x) = x^6 + x^5 + 1$ over $\mathbb{F}_2$ to generate the finite field $\mathbb{F}_{2^6}$. Thus, $\gamma = x$ is a primitive element in $\mathbb{F}_{2^6}$. Let $\beta_i = \gamma^i$ for $1 \leqslant i \leqslant 3$, which is a basis of $\mathbb{F}_{2^6}$ over $\mathbb{F}_4$. Then the generator matrix of the $[9, 3]_{2^6}$ linearized Reed-Solomon code can be given as

$$D = \begin{pmatrix} 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \\ 4 & 8 & 12 & 5 & 9 & 13 & 6 & 10 & 14 \\ 16 & 32 & 48 & 21 & 37 & 53 & 26 & 42 & 58 \end{pmatrix},$$

where the integer $i$ in the matrix stands for the element $\gamma^i \in \mathbb{F}_{2^6}$. Let $\mathcal{C}^*$ be the $[3, 2, 2]_4$ MDS code with generator matrix

$$\begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \in \mathbb{F}_4^{2 \times 3}.$$

By Construction B, the matrix $A_i$ for $1 \leqslant i \leqslant 3$ can be given as

$$A_i = \begin{pmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \end{pmatrix} \in \mathbb{F}_4^{3 \times 5}.$$

Then the generator matrix of the locally repairable codes with information $(2, 2, 3)$-locality can be given as

$$G = (\mathbf{g}_1, \mathbf{g}_2, \ldots, \mathbf{g}_{15})$$
$$= D \operatorname{Diag}(A_1, A_2, A_3)$$
$$= \begin{pmatrix} 1 & 2 & 3 & 59 & 54 & 1 & 2 & 3 & 59 & 54 & 1 & 2 & 3 & 59 & 54 \\ 4 & 8 & 12 & 47 & 27 & 5 & 9 & 13 & 48 & 28 & 6 & 10 & 14 & 49 & 29 \\ 16 & 32 & 48 & 62 & 45 & 21 & 37 & 53 & 4 & 50 & 26 & 42 & 58 & 9 & 55 \end{pmatrix},$$

where the integer $i$ in the matrix stands for the element $\gamma^i \in \mathbb{F}_{2^6}$. Since $\mathrm{Rank}((\mathbf{g}_1, \mathbf{g}_6, \mathbf{g}_{11})) = 3$, we can regard them as information symbols. Their repair sets can be listed as $R_i^j = \{\mathbf{g}_{j+i}, \mathbf{g}_{j+i+2}\}$ for $j \in \{1, 6, 11\}$ and $i = 1, 2$. A computer program verified that indeed the weight of the codewords generated by $G$ is at least 12, i.e., $d = 12 = n - k + 2 - \left\lceil \frac{(k-1)(\delta-1)+1}{(r-1)(\delta-1)+1} \right\rceil$. Thus, the code $\mathcal{C}$ generated by $G$ is a $[15, 3, 12]_{2^6}$ optimal locally repairable codes with

information $(2, 2, 3)$-locality, which is consistent with the result in Corollary 7.

Finally, if $N = 1$, then Construction B is a special case of the construction introduced in [17] for locally repairable codes. In [17], universal and dynamic locally repairable codes with a single repair set and maximal recoverability were considered. In contrast, in Construction B, we mainly focus on locally repairable codes with multiple repair sets.

## VII. CONCLUDING REMARKS

In this paper, a general definition of locality was given that ensures a code symbol can be locally repaired when the number of erasures is bounded by $\delta - 1$. The new definition contains the definitions in [9], [20], [28] as extremal cases. Additionally, a Singleton-type bound was derived for the new codes. Finally, optimal constructions were proposed with respect to the new bound. The constructions can also generate optimal locally repairable codes with information $(r, \delta)_c$-locality, i.e., $(r, \delta - 1, \delta)$-locality with respect to the bound in [28].

However, the codes constructed in this paper have two main drawbacks, namely, low code rates (depending on the number of disjoint repair sets) and large underlying finite fields. One problem that is still open is whether the new bound (like the one in [28]) is also not tight for the high code rate case as shown in [27]. If that is the case, two open questions that remain are how to derive a sharper bound for the high code rate case and how to construct corresponding optimal locally repairable codes. For the low code rate case, the bound in [28] and the new one are tight, but all the known results for those codes require large finite fields. It is very interesting to construct optimal codes with multiple disjoint repair sets over small finite fields, say, of size $O(n)$, as the one proposed in [26] for the single repair set case.

## ACKNOWLEDGMENT

## REFERENCES

[1] A. Agarwal, A. Barg, S. Hu, A. Mazumdar, and I. Tamo, "Combinatorial alphabet-dependent bounds for locally recoverable codes," *IEEE Trans. Inf. Theory*, vol. 64, no. 5, 3481–3492, May 2018.

[2] V. R. Cadambe, C. Huang, and J. Li, "Permutation code: Optimal exact-repair of a single failed node in MDS code based distributed storage systems," in *Proc. IEEE Int. Symp. Inf. Theory*, St. Petersburg, Russia, Jul./Aug. 2011, pp. 1225–1229.

[3] V. R. Cadambe and A. Mazumdar, "Bounds on the size of locally recoverable codes," *IEEE Trans. Inf. Theory*, vol. 61, no. 11, pp. 5787–5794, Nov. 2015.

[4] H. Cai, M. Cheng, C. Fan, and X. Tang, "Optimal locally repairable systematic codes based on packings," *IEEE Trans. Commun.*, vol. 67, no. 1, pp. 39–49, Jan. 2019.

[5] M. Forbes and S. Yekhanin, "On the locality of codeword symbols in non-linear codes," *Discrete Math.*, vol. 324, no. 6, pp. 78–84, Jun. 2014.

[6] C. Fragouli and E. Soljanin, "Network coding fundamentals," *Found. Trends Netw.*, vol. 2, no. 1, pp. 1–133, 2007.

[7] È. M. Gabidulin, "Theory of codes with maximum rank distance," *Problemy Peredachi Informatsii*, vol. 21, no. 1, pp. 3–16, 1985.

[8] P. Gopalan, C. Huang, B. Jenkins, and S. Yekhanin, "Explicit maximally recoverable codes with locality," *IEEE Trans. Inf. Theory*, vol. 60, no. 9, pp. 5245–5256, Sep. 2014.

[9] P. Gopalan, C. Huang, H. Simitci, and S. Yekhanin, "On the locality of codeword symbols," *IEEE Trans. Inf. Theory*, vol. 58, no. 11, pp. 6925–6934, Aug. 2012.

[10] J. Hao and S.-T. Xia, "Constructions of optimal binary locally repairable codes with multiple repair groups," *IEEE Commun. Lett.*, vol. 20, no. 6, pp. 1060–1063, Jun. 2016.

[11] C. Huang, M. Chen, and J. Li, "Pyramid codes: Flexible schemes to trade space for access efficiency in reliable data storage systems," *ACM Trans. Storage*, vol. 9, no. 1, p. 3, 2007.

[12] C. Huang *et al.*, "Erasure coding in windows azure storage," in *Proc. USENIX Assoc.*, 2012, pp. 15–26.

[13] G. Joshi, Y. Liu, and E. Soljanin, "On the delay-storage trade-off in content download from coded distributed storage systems," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 5, pp. 989–997, May 2014.

[14] R. Lidl and H. Niederreiter, *Finite Fields*. Cambridge, U.K.: Cambridge Univ. Press, 1974.

[15] F. J. MacWilliams and N. J. A. Sloane, *The Theory of Error-Correcting Codes*. Amsterdam, The Netherlands: North Holland, 1977.

[16] U. Martínez-Peñas, "Skew and linearized Reed–Solomon codes and maximum sum rank distance codes over any division ring," *J. Algebra*, vol. 504, pp. 587–612, Jun. 2018.

[17] U. Martínez-Peñas and F. R. Kschischang, "Universal and dynamic locally repairable codes with maximal recoverability via sum-rank codes," 2018, *arXiv:1809.11158*. [Online]. Available: https://arxiv.org/abs/1809.11158

[18] L. Pamies-Juarez, H. D. L. Hollmann, and F. Oggier, "Locally repairable codes with multiple repair alternatives," in *Proc. IEEE Int. Symp. Inform. Theory*, Istanbul, Turkey, Jul. 2013, pp. 892–896.

[19] D. S. Papailiopoulos and A. G. Dimakis, "Locally repairable codes," in *Proc. IEEE Int. Symp. Inform. Theory*, Cambridge MA, USA, Jul. 2012, pp. 2771–2775.

[20] N. Prakash, G. M. Kamath, V. Lalitha, and P. V. Kumar, "Optimal linear codes with a local-error-correction property," in *Proc. IEEE Int. Symp. Inform. Theory*, Cambridge MA, USA, Jul. 2012, pp. 2776–2780.

[21] A. S. Rawat, O. O. Koyluoglu, N. Silberstein, and S. Vishwanath, "Optimal locally repairable and secure codes for distributed storage systems," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 212–236, Jan. 2014.

[22] A. S. Rawat, D. S. Papailiopoulos, A. G. Dimakis, and S. Vishwanath, "Locality and availability in distributed storage," *IEEE Trans. Inf. Theory*, vol. 62, no. 8, pp. 4481–4493, Aug. 2016.

[23] N. Silberstein, T. Etzion, and M. Schwartz, "Locality and availability of array codes constructed from subspaces," *IEEE Trans. Inf. Theory*, vol. 65, no. 5, pp. 2648–2660, May 2019.

[24] W. Song, S. H. Dau, C. Yuen, and T. J. Li, "Optimal locally repairable linear codes," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 5, pp. 1019–1036, May 2014.

[25] Y. S. Su, "On the construction of local parties for $(r, t)$-availability in distributed storage," *IEEE Trans. Commun.*, vol. 65, no. 6, pp. 2332–2344, Jun. 2017.

[26] I. Tamo and A. Barg, "A family of optimal locally recoverable codes," *IEEE Trans. Inf. Theory*, vol. 60, no. 8, pp. 4661–4676, Aug. 2014.

[27] I. Tamo, A. Barg, and A. Frolov, "Bounds on the parameters of locally recoverable codes," *IEEE Trans. Inf. Theory*, vol. 62, no. 6, pp. 3070–3083, Jun. 2016.

[28] A. Wang and Z. Zhang, "Repair locality with multiple erasure tolerance," *IEEE Trans. Inf. Theory*, vol. 60, no. 11, pp. 6979–6987, Nov. 2014.

[29] A. Wang and Z. Zhang, "An integer programming-based bound for locally repairable codes," *IEEE Trans. Inf. Theory*, vol. 61, no. 10, pp. 5280–5294, Oct. 2015.

**Han Cai** (S'16–M'18) received the B.S. and M.S. degrees in mathematics from Hubei University, Wuhan, China, in 2009 and 2013, respectively and received the Ph.D. degree from the department of communication engineering, Southwest Jiaotong University, Chengdu, China, in 2017. During Oct. 2015 to Oct. 2017, he was a visiting Ph.D. student in the Faculty of Engineering, Information and Systems, University of Tsukuba, Japan. He is currently a postdoctoral fellow at the Department of Electrical & Computer Engineering, Ben-Gurion University of the Negev, Israel. His research interests include coding theory and sequence design.

**Ying Miao** received the D.Sci. degree in mathematics from Hiroshima University, Hiroshima, Japan, in 1997. From 1989 to 1993, he worked for Suzhou Institute of Silk Textile Technology, Suzhou, Jiangsu, P. R. China. From 1995 to 1997, he was a Research Fellow of the Japan Society for the Promotion of Science. During 1997–1998, he was a Postdoctoral Fellow at the Department of Computer Science, Concordia University, Montreal, QC, Canada. In 1998, he joined the University of Tsukuba, Tsukuba, Ibaraki, Japan, where he is currently a Full Professor at the Faculty of Engineering, Information and Systems. His current research interests include combinatorics, coding theory, and information security.

Dr. Miao is on the Editorial Boards of several journals such as *Graphs and Combinatorics*, and *Journal of Combinatorial Designs*. He received the 2001 Kirkman Medal from the Institute of Combinatorics and its Applications.

**Moshe Schwartz** (M'03–SM'10) is an associate professor at the School of Electrical and Computer Engineering, Ben-Gurion University of the Negev, Israel. His research interests include algebraic coding, combinatorial structures, and digital sequences.

Prof. Schwartz received the B.A. *(summa cum laude)*, M.Sc., and Ph.D. degrees from the Technion–Israel Institute of Technology, Haifa, Israel, in 1997, 1998, and 2004 respectively, all from the Computer Science Department. He was a Fulbright post-doctoral researcher in the Department of Electrical and Computer Engineering, University of California San Diego, and a post-doctoral researcher in the Department of Electrical Engineering, California Institute of Technology. While on sabbatical 2012–2014, he was a visiting scientist at the Massachusetts Institute of Technology (MIT).

Prof. Schwartz received the 2009 IEEE Communications Society Best Paper Award in Signal Processing and Coding for Data Storage. He has also been serving as an Associate Editor for Coding Techniques for the IEEE TRANSACTIONS ON INFORMATION THEORY since 2014.

**Xiaohu Tang** (M'04–SM'18) received the B.S. degree in applied mathematics from the Northwest Polytechnic University, Xi'an, China, the M.S. degree in applied mathematics from the Sichuan University, Chengdu, China, and the Ph.D. degree in electronic engineering from the Southwest Jiaotong University, Chengdu, China, in 1992, 1995, and 2001 respectively.

From 2003 to 2004, he was a research associate in the Department of Electrical and Electronic Engineering, Hong Kong University of Science and Technology. From 2007 to 2008, he was a visiting professor at University of Ulm, Germany. Since 2001, he has been in the School of Information Science and Technology, Southwest Jiaotong University, where he is currently a professor. His research interests include coding theory, network security, distributed storage and information processing for big data.

Dr. Tang was the recipient of the National excellent Doctoral Dissertation award in 2003 (China), the Humboldt Research Fellowship in 2007 (Germany), and the Outstanding Young Scientist Award by NSFC in 2013 (China). He served as Associate Editors for several journals including IEEE TRANSACTIONS ON INFORMATION THEORY and *IEICE Transactions on Fundamentals*, and served on a number of technical program committees of conferences.