# Capacity of Private Linear Computation for Coded Databases

Sarah A. Obead[†], Hsuan-Yin Lin[‡], Eirik Rosnes[‡], and Jörg Kliewer[†]
[†]Helen and John C. Hartmann Department of Electrical and Computer Engineering
New Jersey Institute of Technology, Newark, New Jersey 07102, USA
[‡]Simula UiB, N–5020 Bergen, Norway

*Abstract*—We consider the problem of private linear computation (PLC) in a distributed storage system. In PLC, a user wishes to compute a linear combination of $f$ messages stored in noncolluding databases while revealing no information about the coefficients of the desired linear combination to the databases. In extension of our previous work we employ linear codes to encode the information on the databases. We show that the PLC capacity, which is the ratio of the desired linear function size and the total amount of downloaded information, matches the maximum distance separable (MDS) coded capacity of private information retrieval for a large class of linear codes that includes MDS codes. In particular, the proposed converse is valid for any number of messages and linear combinations, and the capacity expression depends on the rank of the coefficient matrix obtained from all linear combinations.

## I. INTRODUCTION

Private information retrieval (PIR) from public databases has been the focus of attention for several decades in the computer science community (see, e.g., [1], [2]). Recently, the aspect of minimizing the communication cost, e.g., the required rate or bandwidth of privately querying the databases with the desired requests and downloading the corresponding information from the databases has in particular been of interest also in the information theory and coding communities (see, e.g., [3]–[10]). Those settings primarily focus on distributed storage systems (DSSs), where data is encoded by a linear code and then stored across several databases.

A recently proposed generalization of the PIR problem [11]–[15] addresses the *private computation (PC)* of functions of these messages. In PC a user has access to a given number of databases and intends to compute a function of messages stored in these databases. This function is kept private from the databases, as they may be under the control of an adversary. In [11], [12], the authors consider that the data is replicated at each database, which are assumed to be noncolluding. Both works characterize the fundamental information-theoretic communication overhead needed to privately compute a given *linear* function, called private linear computation (PLC), and specify the corresponding capacity and achievable rates as a function of the number

of messages and the number of databases, respectively. Interestingly, the obtained PLC capacity is equal to the PIR capacity of [6]. The extension to the coded case is addressed in [13], [14]. In [13], a PC scheme for polynomial functions of the data over $t$ colluding databases is constructed by generalizing the star-product PIR scheme of [9]. Specifically, for systematically coded storage, this scheme considers functions that are polynomials of fixed degree $g$, and achieves a rate matching the best asymptotic rate (when the number of messages goes to infinity) of coded PIR $R = 1 - \gamma_c$ for $g = t = 1$ (the case of linear function retrieval and noncolluding databases), where $\gamma_c$ denotes the code rate. In our previous work [14], we consider maximum distance separable (MDS) coded storage and the computation of linear message functions. The presented scheme is able to achieve the MDS-coded PIR capacity in [8], referred to as the MDS-PIR capacity in the sequel, and strictly generalizes the achievable schemes in [11], [12].

On the other hand, for the case of noncolluding databases, a PIR protocol for a DSS where data is stored using a non-MDS linear code, is proposed in [10]. The protocol is shown to achieve the asymptotic MDS-PIR capacity for several example codes. In [16], [17], this approach is extended to colluding databases and it is proved that both the asymptotic and the nonasymptotic MDS-PIR capacity can be achieved for a large class of linear codes in the noncolluding case. For noncolluding databases, the first family of non-MDS codes for which the PIR capacity is known is found in [18], [19]. Further, PIR on linearly-coded databases for the case of colluding databases is also addressed in [5], [9], [20]. For the PLC case however, capacity results for arbitrary linear codes have not been addressed so far in the open literature to the best of our knowledge.

In this work, we intend to fill this void and propose an extension of our previous works in [14] and [17] to PLC for arbitrary linear codes with noncolluding databases. Surprisingly, we show that the PLC capacity for a large class of linear codes matches the MDS-PIR capacity when a user wishes to compute an arbitrary linear combination of $f$ independent equal-sized messages over some finite field $\mathbb{F}_q$, distributed over $n$ noncolluding linearly-coded databases. In particular, by extending our previous work [14] where only achievable rates are stated, we also adapt the converse proof

of [19, Thm. 4] to the PLC problem in this paper. In contrast to [12], our converse is valid for any number of messages $f$ and any number of desired linear combinations $\mu$, and our capacity result depends on the rank of the coefficient matrix obtained from all $\mu$ linear combinations.

## II. DEFINITIONS AND PROBLEM STATEMENT

### A. Notation and Definitions

We denote by $\mathbb{N}$ the set of all positive integers and for some $a, b \in \mathbb{N}$, $a \leq b$, $[a : b] \triangleq \{a, a+1, \ldots, b\}$. Random and deterministic quantities are carefully distinguished as follows. A random variable is denoted by a capital Roman letter, e.g., $X$, while its realization is denoted by the corresponding small Roman letter, e.g., $x$; vectors are boldfaced, e.g., $\boldsymbol{X}$ denotes a random vector and $\boldsymbol{x}$ denotes a deterministic vector; random matrices are represented by bold sans serif letters, e.g., $\mathsf{X}$, and $\mathsf{X}$ represents its realization. In addition, sets are denoted by calligraphic upper case letters, e.g., $\mathcal{X}$, and $\mathcal{X}^c$ denotes the complement of a set $\mathcal{X}$ in a universe set. For a given index set $\mathcal{S}$, we also write $\mathsf{X}^{\mathcal{S}}$ and $Y_{\mathcal{S}}$ to represent $\{\mathsf{X}^{(v)} : v \in \mathcal{S}\}$ and $\{Y_j : j \in \mathcal{S}\}$, respectively. Furthermore, some constants are also depicted by Greek letters or a special font, e.g., $\mathsf{X}$. We denote a submatrix of $\mathsf{X}$ that is restricted in columns by the set $\mathcal{I}$ by $\mathsf{X}|_{\mathcal{I}}$. The function $\mathsf{H}(\cdot)$ represents the entropy of its argument and $\mathsf{I}(X; Y)$ denotes the mutual information of the random variables $X$, $Y$. $(\cdot)^{\mathsf{T}}$ denotes the transpose of its argument, while $\mathrm{rank}(\mathsf{V})$ denotes the rank of a matrix $\mathsf{V}$. We use the customary code parameters $[n, k]$ to denote a code $\mathscr{C}$ over the finite field $\mathbb{F}_q$ of blocklength $n$ and dimension $k$. A generator matrix of $\mathscr{C}$ is denoted by $\mathsf{G}^{\mathscr{C}}$. A set of coordinates of $\mathscr{C}$, $\mathcal{I} \subseteq [1 : n]$, of size $k$ is said to be an *information set* if and only if $\mathsf{G}^{\mathscr{C}}|_{\mathcal{I}}$ is invertible. The function $\chi(\boldsymbol{x})$ denotes the support of a vector $\boldsymbol{x}$.

### B. System Model

The PLC problem for coded DSSs is described as follows. We consider a DSS that stores $f$ messages $\mathsf{W}^{(1)}, \ldots, \mathsf{W}^{(f)}$, where each message $\mathsf{W}^{(m)} = (W_{i,j}^{(m)})$, $m \in [1 : f]$, can be seen as a random matrix of size $\beta \times k$ over $\mathbb{F}_q$ with $\beta, k \in \mathbb{N}$. Each message is encoded using an $[n, k]$ code as follows. Let $\boldsymbol{W}_i^{(m)} = (W_{i,1}^{(m)}, \ldots, W_{i,k}^{(m)})$, $i \in [1 : \beta]$, be a message vector corresponding to the $i$-th row of $\mathsf{W}^{(m)}$. Each $\boldsymbol{W}_i^{(m)}$ is encoded by an $[n, k]$ code $\mathscr{C}$ over $\mathbb{F}_q$ into a length-$n$ codeword $\boldsymbol{C}_i^{(m)} = (C_{i,1}^{(m)}, \ldots, C_{i,n}^{(m)})$. The $\beta f$ generated codewords $\boldsymbol{C}_i^{(m)}$ are then arranged in the array $\mathsf{C} = ((\mathsf{C}^{(1)})^{\mathsf{T}} | \ldots | (\mathsf{C}^{(f)})^{\mathsf{T}})^{\mathsf{T}}$ of dimensions $\beta f \times n$, where $\mathsf{C}^{(m)} = ((\boldsymbol{C}_1^{(m)})^{\mathsf{T}} | \ldots | (\boldsymbol{C}_{\beta}^{(m)})^{\mathsf{T}})^{\mathsf{T}}$. The code symbols $C_{1,j}^{(m)}, \ldots, C_{\beta,j}^{(m)}$, $m \in [1 : f]$, for all $f$ messages are stored on the $j$-th database, $j \in [1 : n]$.

Let $\mathsf{L} \triangleq \beta \cdot k$. Then, each message $\mathsf{W}^{(m)}$, $m \in [1 : f]$, can also be seen as a random vector variable $\mathsf{W}^{(m)} = (W_1^{(m)}, \ldots, W_{\mathsf{L}}^{(m)})$ of $\mathsf{L}$ symbols that are chosen independently and uniformly at random from $\mathbb{F}_q$. Thus,

$$\mathsf{H}(\mathsf{W}^{(m)}) = \mathsf{L}, \; \forall \, m \in [1 : f],$$

$$\mathsf{H}(\mathsf{W}^{(1)}, ..., \mathsf{W}^{(f)}) = f\mathsf{L} \quad \text{(in } q\text{-ary units)}.$$

### C. Private Linear Computation for Coded DSSs

We consider the case where no set of databases can collude. A user wishes to privately compute exactly one linear function from the $\mu$ *candidate* linear functions $\mathsf{X}^{(1)}, \ldots, \mathsf{X}^{(\mu)}$ from the coded DSS, where $\mathsf{X}^{(v)} = (X_1^{(v)}, \ldots, X_{\mathsf{L}}^{(v)})$. The $\mu$-tuple $(X_l^{(1)}, \ldots, X_l^{(\mu)})^{\mathsf{T}}$, $\forall \, l \in [1 : \mathsf{L}]$, is mapped by a deterministic matrix $\mathsf{V}$ of size $\mu \times f$ over $\mathbb{F}_q$ by

$$\begin{pmatrix} X_l^{(1)} \\ \vdots \\ X_l^{(\mu)} \end{pmatrix} = \mathsf{V}_{\mu \times f} \begin{pmatrix} W_l^{(1)} \\ \vdots \\ W_l^{(f)} \end{pmatrix}. \tag{1}$$

Hence, the user privately generates an index $v \in [1 : \mu]$ and wishes to compute the $v$-th linear function while keeping the index $v$ private from each database. Here, we also assume that the rank of $\mathsf{V}$ is equal to $\mathrm{rank}(\mathsf{V}) = r \leq \min\{\mu, f\}$ and the indices corresponding to a basis for the row space of $\mathsf{V}$ are denoted by the set $\mathcal{L} \triangleq \{\ell_1, \ldots, \ell_r\} \subseteq [1 : \mu]$.

In order to retrieve the linear function $\mathsf{X}^{(v)}$ from the coded DSS, $v \in [1 : \mu]$, the user sends a query $Q_j^{(v)}$ over $\mathbb{F}_q$ to the $j$-th database for all $j \in [1 : n]$. Since the queries are generated by the user without any prior knowledge of the realizations of the candidate linear functions, the queries are independent of the linear functions. In other words, we have

$$\mathsf{I}(\mathsf{X}^{(1)}, \ldots, \mathsf{X}^{(\mu)}; Q_1^{(v)}, \ldots, Q_n^{(v)}) = 0, \, \forall \, v \in [1 : \mu].$$

In response to the received query, database $j$ sends the answer $A_j^{(v)}$ back to the user. $A_j^{(v)}$ is a deterministic function of $Q_j^{(v)}$ and the data stored in the database. Thus, $\forall \, v \in [1 : \mu]$,

$$\mathsf{H}(A_j^{(v)} \mid Q_j^{(v)}, \mathbf{C}^{[1:f]}) = 0, \, \forall \, j \in [1 : n].$$

Next, we define a PLC protocol for coded DSSs as follows.

*Definition 1:* Consider a DSS with $n$ noncolluding databases storing $f$ messages using an $[n, k]$ code. The user wishes to retrieve the $v$-th linear function $\mathsf{X}^{(v)}$, $v \in [1 : \mu]$, from the available information $Q_j^{(v)}$ and $A_j^{(v)}$, $j \in [1 : n]$. For a PLC protocol, the following conditions must be satisfied $\forall \, v \in [1 : \mu]$,

[Privacy]
$$\mathsf{I}(v; Q_j^{(v)}, A_j^{(v)}, \mathsf{X}^{(1)}, \ldots, \mathsf{X}^{(\mu)}) = 0, \, \forall \, j \in [1 : n],$$

[Recovery]
$$\mathsf{H}(\mathsf{X}^{(v)} \mid A_1^{(v)}, \ldots, A_n^{(v)}, Q_1^{(v)}, \ldots, Q_n^{(v)}) = 0.$$

To measure the efficiency of a PLC protocol, we consider the required number of downloaded symbols for retrieving the $\mathsf{L}$ symbols of the candidate linear function.

*Definition 2 (PLC Rate and Capacity for Coded DSSs):* The rate of a PLC protocol, denoted by $\mathsf{R}$, is the ratio of the desired linear function size $\mathsf{L}$ to the total required download cost $\mathsf{D}$, i.e.,

$$\mathsf{R} \triangleq \frac{\mathsf{L}}{\mathsf{D}}.$$

The PLC *capacity* $\mathsf{C}_{\mathrm{PLC}}$ is the maximum possible PLC rate over all possible PLC protocols for a given $[n, k]$ storage code.

### D. MDS-PIR Capacity-Achieving Codes

In [16], [17], a PIR protocol for any linearly-coded DSS that uses an $[n,k]$ code to store $f$ messages, named Protocol 1, is proposed. The corresponding PIR rate depends on the following property of the underlying storage code $\mathscr{C}$.

*Definition 3:* Let $\mathscr{C}$ be an arbitrary $[n,k]$ code. A $\nu \times n$ binary matrix $\Lambda_{\kappa,\nu}(\mathscr{C})$ is said to be a *PIR achievable rate matrix* for $\mathscr{C}$ if the following conditions are satisfied.

1) The Hamming weight of each column of $\Lambda_{\kappa,\nu}$ is $\kappa$, and
2) for each matrix row $\boldsymbol{\lambda}_i$, $i \in [1:\nu]$, $\chi(\boldsymbol{\lambda}_i)$ always contains an information set.

In other words, each coordinate $j$ of $\mathscr{C}$, $j \in [1:n]$, appears exactly $\kappa$ times in $\{\chi(\boldsymbol{\lambda}_i)\}_{i \in [1:\nu]}$, and every set $\chi(\boldsymbol{\lambda}_i)$ contains an information set.

In [16], [17], it is shown that the MDS-PIR capacity (i.e., the PIR capacity for a DSS where data is encoded and stored using an MDS code) can be achieved using Protocol 1 for a special class of $[n,k]$ codes. In particular, to achieve the MDS-PIR capacity using Protocol 1, the $[n,k]$ storage code should possess a specific underlying structure as given by the following theorem.

*Theorem 1 ([16], [17]):* Consider a DSS that uses an $[n,k]$ code $\mathscr{C}$ to store $f$ messages. If a PIR achievable rate matrix $\Lambda_{\kappa,\nu}(\mathscr{C})$ with $\frac{\kappa}{\nu} = \frac{k}{n}$ exists, then the MDS-PIR capacity

$$\mathsf{C}_{\text{MDS-PIR}} \triangleq \left(1 - \frac{k}{n}\right)\left[1 - \left(\frac{k}{n}\right)^f\right]^{-1}$$

is achievable.

This gives rise to the following definition.

*Definition 4:* A PIR achievable rate matrix $\Lambda_{\kappa,\nu}(\mathscr{C})$ with $\frac{\kappa}{\nu} = \frac{k}{n}$ for an $[n,k]$ code $\mathscr{C}$ is called an *MDS-PIR capacity-achieving* matrix, and $\mathscr{C}$ is referred to as an *MDS-PIR capacity-achieving* code.

In Section IV, we will present a PLC protocol and a general achievable rate for any $\text{rank}(\mathsf{V}) = r$ by using the PIR achievable rate matrix $\Lambda_{\kappa,\nu}$ of an $[n,k]$ code.

### III. PLC CAPACITY FOR MDS-PIR CAPACITY-ACHIEVING CODES

The main result of this paper is the derivation of the PLC capacity for a coded DSS where data is encoded and stored using a linear code from the class of MDS-PIR capacity-achieving codes introduced in [16], [17], which is stated in the following theorem.

*Theorem 2:* Consider a DSS that uses an $[n,k]$ MDS-PIR capacity-achieving code $\mathscr{C}$ to store $f$ messages. Then, the maximum achievable PLC rate over all possible PLC protocols, i.e., the PLC capacity, is

$$\mathsf{C}_{\text{PLC}} \triangleq \left(1 - \frac{k}{n}\right)\left[1 - \left(\frac{k}{n}\right)^r\right]^{-1} \qquad (2)$$

for any rank $r$ of the linear mapping from (1).

The achievable scheme is provided in Section IV, while the converse proof is given in Section V. Note that the class of MDS-PIR capacity-achieving codes includes MDS

codes, cyclic codes, Reed-Muller codes, and certain classes of distance-optimal local reconstruction codes [16], [17].

We remark here that it is known that if $\text{rank}(\mathsf{V}) = f$, then the PLC capacity for an MDS-coded DSS is equal to the MDS-PIR capacity [14], i.e., if $\text{rank}(\mathsf{V}) = f$, then

$$\mathsf{C}_{\text{PLC}} = \left(1 - \frac{k}{n}\right)\left[1 - \left(\frac{k}{n}\right)^f\right]^{-1} = \mathsf{C}_{\text{MDS-PIR}}.$$

### IV. ACHIEVABLE SCHEME

We start with a coded PIR query scheme for a message of $\mu$ *dependent* virtual messages. Given that the messages are stored using an $[n,k]$ MDS-PIR capacity-achieving code $\mathscr{C}$, we find the associated $\nu \times n$ MDS-PIR capacity-achieving matrix $\Lambda_{\kappa,\nu}$ and obtain the PIR interference matrices $\mathsf{A}_{\kappa \times n}$ and $\mathsf{B}_{(\nu-\kappa) \times n}$ as given by the following definition.

*Definition 5 ([16], [17]):* For a given $\nu \times n$ PIR achievable rate matrix $\Lambda_{\kappa,\nu}(\mathscr{C}) = (\lambda_{u,j})$, we define the PIR interference matrices $\mathsf{A}_{\kappa \times n} = (a_{i,j})$ and $\mathsf{B}_{(\nu-\kappa) \times n} = (b_{i,j})$ for the code $\mathscr{C}$ with

$$a_{i,j} \triangleq u \text{ if } \lambda_{u,j} = 1, \forall j \in [1:n], i \in [1:\kappa], u \in [1:\nu],$$
$$b_{i,j} \triangleq u \text{ if } \lambda_{u,j} = 0, \forall j \in [1:n], i \in [1:\nu-\kappa], u \in [1:\nu].$$

Note that in Definition 5, for each $j \in [1:n]$, distinct values of $u \in [1:\nu]$ should be assigned for all $i$. Thus, the assignment is not unique in the sense that the order of the entries of each column of A and B can be permuted. Next, for the sake of illustrating our query set generation algorithm, we make use of the following definition.

*Definition 6:* By $\mathcal{S}(u|\mathsf{A}_{\kappa \times n})$ we denote the set of column coordinates of matrix $\mathsf{A}_{\kappa \times n} = (a_{i,j})$ in which at least one of its entries is equal to $u$, i.e.,

$$\mathcal{S}(u|\mathsf{A}_{\kappa \times n}) \triangleq \{j \in [1:n] : \exists a_{i,j} = u, i \in [1:\kappa]\}.$$

As a result, we require the size of the message to be $\mathsf{L} = \nu^\mu \cdot k$ (i.e., $\beta = \nu^\mu$). The scheme is completed in $\mu$ rounds.

### A. Query Generation

Before running the main algorithm to generate query sets, the following index preparation for the coded symbols stored in each database is performed.

***1) Index Preparation:*** The goal is to make the coded symbols queried from each database to appear to be chosen randomly and independently from the desired linear function index. Note that the linear function is computed separately for the $t$-th row of all messages, $t \in [1:\beta]$. Therefore, similar to the PLC scheme in [12] and the MDS-coded PLC scheme in [14], we apply a permutation that is fixed across all coded symbols for the $t$-th row to maintain the dependency across the associated message elements. Let $\pi(\cdot)$ be a random permutation function over $[1:\beta]$, and let

$$U_{t,j}^{(v)} \triangleq \boldsymbol{v}_v \boldsymbol{C}_{\pi(t),j}, \quad t \in [1:\beta], \quad j \in [1:n],$$

denote the $t$-th permuted code symbol associated with the $v$-th virtual message $\mathbf{X}^{(v)}$ stored in the $j$-th database. Here, $\boldsymbol{v}_v$ represents the $v$-th row vector of the matrix $\mathsf{V}_{\mu \times f} = (v_{i,j})$ and $\boldsymbol{C}_{t,j} \triangleq \left(C_{t,j}^{(1)}, \ldots, C_{t,j}^{(f)}\right)^\intercal$. The permutation $\pi(\cdot)$ is randomly selected privately and uniformly by the user.

*2) Preliminaries*: For $\tau \in [1:\mu]$, a sum $U_{i_1,j}^{(v_1)} + U_{i_2,j}^{(v_2)} + \cdots + U_{i_\tau,j}^{(v_\tau)}$, $j \in [1:n]$, of $\tau$ distinct symbols is a $\tau$-sum for any subset $\{i_1, i_2, \ldots, i_\tau\} \subseteq [1:\beta]$, and $\{v_1, v_2, \ldots, v_\tau\} \subseteq [1:\mu]$ determines the *type* of the $\tau$-sum. Since we have $\binom{\mu}{\tau}$ different selections of $\tau$ distinct elements out of $\mu$ elements, a $\tau$-sum can have $\binom{\mu}{\tau}$ different types. The query generation procedure is subdivided into $\mu$ rounds. For a requested linear function indexed by $v \in [1:\mu]$, a query set $Q_j^{(v)}$, $j \in [1:n]$, is composed of $\mu$ disjoint subsets, one generated by each round $\tau \in [1:\mu]$. For each round the query subset is further subdivided into two subsets. The first subset $Q_j^{(v)}(\mathcal{D}; \tau)$ consists of $\tau$-sums with a single symbol from the *desired* linear function and $\tau - 1$ symbols from *undesired* linear functions, while the second subset $Q_j^{(v)}(\mathcal{U}; \tau)$ contains $\tau$-sums with symbols only from undesired linear functions. To this end, in the scheme we introduce $\kappa n$ auxiliary query sets $Q_j^{(v)}(a_{i,j}, \mathcal{D}; \tau)$, $i \in [1:\kappa]$, where each query consists of a distinct symbol from the desired linear function and $\tau - 1$ symbols from undesired linear functions, and $(\nu - \kappa)n$ auxiliary query sets $Q_j^{(v)}(b_{i,j}, \mathcal{U}; \tau)$, $i \in [1:\nu - \kappa]$, to represent the side information query sets for each database $j \in [1:n]$. We utilize these sets to generate the query sets of each round according to the PIR interference matrices $\mathsf{A}_{\kappa \times n}$ and $\mathsf{B}_{(\nu - \kappa) \times n}$. The query sets for all databases are generated by Algorithm 1 through the following procedures.[1]

*3) Initialization (Round $\tau = 1$)*: In the initialization step, the algorithm generates the auxiliary queries for the first round. This round is described by Steps 5 to 11 of Algorithm 1, where we have $\tau = 1$ for the $\tau$-sum. At this point, Algorithm 1 invokes the subroutine `Initial-Round` given in Algorithm 2 to generate $Q_j^{(v)}(a_{i,j}, \mathcal{D}; 1)$, $i \in [1:\kappa]$, such that each of these query sets contains $\alpha_1 = \kappa^{\mu - 1}$ distinct code symbols. Furthermore, to maintain function symmetry, the algorithm asks each database for the same number of distinct symbols of all other linear functions in $Q_j^{(v)}(a_{i,j}, \mathcal{U}; 1)$, $i \in [1:\kappa]$, resulting in a total number of $\binom{\mu - 1}{1}\kappa^{\mu - 1}$ symbols. As a result, the queried code symbols in the auxiliary query sets for each database are symmetric with respect to all linear function vectors indexed by $v' \in [1:\mu]$. We associate the symbols of undesired functions in $\kappa$ groups, each placed in the undesired query sets $Q_j^{(v)}(a_{i,j}, \mathcal{U}; 1)$, $i \in [1:\kappa]$, to be exploited as distinct side information for the second round of the scheme. Since this procedure produces $\kappa$ undesired query sets for each database, database symmetry is maintained.

*4) Desired Function Symbols for Rounds $\tau > 1$*: For the following rounds a similar process is repeated in terms of generating auxiliary query sets containing distinct code symbols from the desired linear function $\mathbf{U}^{(v)} = (U_{t,j}^{(v)})$. This is accomplished in Steps 16 to 18 by calling the subroutine `Desired-Q`, given in Algorithm 3, to generate

---

[1]Note that a query $Q_j^{(v)}$ sent to the $j$-th database usually indicates the row indices of the stored code symbols that the user requests, while the answer $A_j^{(v)}$ to the query $Q_j^{(v)}$ refers to the particular code symbols requested through the query. In Algorithm 1, with some abuse of notation, the generated queries are sets containing their answers.

---

**Algorithm 1:** `Q-Gen`

---

**Input** : $v$, $\mu$, $\kappa$, $\nu$, $n$, $\mathsf{A}_{\kappa \times n}$, and $\mathsf{B}_{(\nu - \kappa) \times n}$
**Output:** $Q_1^{(v)}, \ldots, Q_n^{(v)}$

1 **for** $\tau \in [1:\mu]$ **do**
2    $Q_j^{(v)}(\mathcal{D}; \tau) \leftarrow \emptyset$, $Q_j^{(v)}(\mathcal{U}; \tau) \leftarrow \emptyset$, $j \in [1:n]$
3    $\alpha_\tau \leftarrow \kappa^{\mu - 1} + \sum_{h=1}^{\tau - 1} \binom{\mu - 1}{h} \kappa^{\mu - (h+1)}(\nu - \kappa)^h$
4    ▷ Generate query sets for the initial round
5    **if** $\tau = 1$ **then**
6      **for** $u \in [1:\nu]$ **do**
7        **for** $j \in \mathcal{S}(u | \mathsf{A}_{\kappa \times n})$ **do**
8          $Q_j^{(v)}(u, \mathcal{D}; \tau), Q_j^{(v)}(u, \mathcal{U}; \tau) \leftarrow$ `Initial-Round`$(u, \alpha_\tau, j, v, \tau)$
9        **end**
10      **end**
11    **end**
12    ▷ Generate query sets for the following rounds $\tau > 1$
13    **else**
14      **for** $u \in [1:\nu]$ **do**
15        ▷ Generate desired symbols for the following rounds $\tau > 1$
16        **for** $j \in \mathcal{S}(u | \mathsf{A}_{\kappa \times n})$ **do**
17          $Q_j^{(v)}(u, \mathcal{D}; \tau) \leftarrow$ `Desired-Q`$(u, \alpha_\tau, j, v, \tau)$
18        **end**
19        ▷ Generate side information for the following rounds $\tau > 1$
20        **for** $j \in \mathcal{S}(u | \mathsf{B}_{(\nu - \kappa) \times n})$ **do**
21          $Q_j^{(v)}(u, \mathcal{U}; \tau - 1) \leftarrow$ `Exploit-SI`$(u, Q_1^{(v)}(u, \mathcal{U}, \tau - 1), \ldots, Q_n^{(v)}(u, \mathcal{U}, \tau - 1), j, v, \tau)$
22        **end**
23      **end**
24      ▷ Generate the final desired query sets for the following rounds $\tau > 1$
25      **for** $j \in [1:n]$ **do**
26        $\tilde{Q}_j^{(v)}(\mathcal{U}; \tau - 1) \leftarrow \bigcup_{i \in [1:\nu - \kappa]} Q_j^{(v)}(b_{i,j}, \mathcal{U}; \tau - 1)$
27        $\tilde{Q}_j^{(v)}(1, \mathcal{U}; \tau - 1), \ldots, \tilde{Q}_j^{(v)}(\kappa, \mathcal{U}; \tau - 1) \leftarrow$ `Partition`$(\tilde{Q}_j^{(v)}(\mathcal{U}; \tau - 1))$
28        **for** $i \in [1:\kappa]$ **do**
29          $Q_j^{(v)}(a_{i,j}, \mathcal{D}; \tau) \leftarrow$ `SetAddition`$(Q_j^{(v)}(a_{i,j}, \mathcal{D}; \tau), \tilde{Q}_j^{(v)}(i, \mathcal{U}; \tau - 1))$
30        **end**
31      **end**
32      ▷ Generate the query sets of undesired symbols by forcing message symmetry for the following rounds $\tau > 1$
33      **for** $u \in [1:\nu]$ **do**
34        **for** $j \in \mathcal{S}(u | \mathsf{A}_{\kappa \times n})$ **do**
35          $Q_j^{(v)}(u, \mathcal{U}; \tau) \leftarrow$ `M-Sym`$(Q_j^{(v)}(u, \mathcal{D}; \tau), j, v, \tau)$
36        **end**
37      **end**
38    **end**
39    **for** $u \in [1:\nu]$ **do**
40      **for** $j \in \mathcal{S}(u | \mathsf{A}_{\kappa \times n})$ **do**
41        $Q_j^{(v)}(\mathcal{D}; \tau) \leftarrow Q_j^{(v)}(\mathcal{D}; \tau) \cup Q_j^{(v)}(u, \mathcal{D}; \tau)$
42        $Q_j^{(v)}(\mathcal{U}; \tau) \leftarrow Q_j^{(v)}(\mathcal{U}; \tau) \cup Q_j^{(v)}(u, \mathcal{U}; \tau)$
43      **end**
44    **end**
45 **end**
46 **for** $j \in [1:n]$ **do**
47    $Q_j^{(v)} \leftarrow \bigcup_{\tau \in [1:\mu]} \left(Q_j^{(v)}(\mathcal{D}; \tau) \cup Q_j^{(v)}(\mathcal{U}; \tau)\right)$
48 **end**

**Algorithm 2: Initial-Round**

> **Input** : $u$, $\alpha_\tau$, $j$, $v$, and $\tau$
> **Output:** $\phi^{(v)}(u, \mathcal{D}; \tau), \phi^{(v)}(u, \mathcal{U}; \tau)$
> 1  $\phi^{(v)}(u, \mathcal{D}; \tau) \leftarrow \emptyset$, $\phi^{(v)}(u, \mathcal{U}; \tau) \leftarrow \emptyset$
> 2  **for** $l \in [1 : \alpha_\tau]$ **do**
> 3  $\quad \phi^{(v)}(u, \mathcal{D}; \tau) \leftarrow \phi^{(v)}(u, \mathcal{D}; \tau) \cup \{U^{(v)}_{(u-1) \cdot \alpha_\tau + l, j}\}$
> 4  $\quad \phi^{(v)}(u, \mathcal{U}; \tau) \leftarrow \phi^{(v)}(u, \mathcal{U}; \tau) \cup$
>    $\qquad \left( \bigcup_{v'=1}^{\mu} \{U^{(v')}_{(u-1) \cdot \alpha_\tau + l, j}\} \setminus \{U^{(v)}_{(u-1) \cdot \alpha_\tau + l, j}\} \right)$
> 5  **end**

$Q^{(v)}_j(a_{i,j}, \mathcal{D}; \tau)$, $i \in [1 : \kappa]$, such that each of these query sets contains $(\alpha_\tau - 1) - \alpha_{\tau-1} + 1 = \binom{\mu-1}{\tau-1} \kappa^{\mu-(\tau-1+1)} (\nu - \kappa)^{\tau-1}$ distinct code symbols from the desired linear function $\mathbf{U}^{(v)}$.

**Algorithm 3: Desired-Q**

> **Input** : $u$, $\alpha_\tau$, $j$, $v$, and $\tau$
> **Output:** $\phi^{(v)}(u, \mathcal{D}; \tau)$
> 1  $\phi^{(v)}(u, \mathcal{D}; \tau) \leftarrow \emptyset$
> 2  **for** $l \in [\alpha_{\tau-1} : \alpha_\tau - 1]$ **do**
> 3  $\quad \phi^{(v)}(u, \mathcal{D}; \tau) \leftarrow \phi^{(v)}(u, \mathcal{D}; \tau) \cup \{U^{(v)}_{l \cdot \nu + u, j}\}$
> 4  **end**

*5) Side Information Exploitation:* In Steps 20 to 22, we generate the side information query sets $Q^{(v)}_j(b_{i',j}, \mathcal{U}; \tau - 1)$, $i' \in [1 : \nu - \kappa]$, from the auxiliary query sets $Q^{(v)}_1(a_{i,1}, \mathcal{U}; \tau - 1), \ldots, Q^{(v)}_n(a_{i,n}, \mathcal{U}; \tau - 1)$, $i \in [1 : \kappa]$, of the previous round $\tau - 1$, $\tau \in [2 : \mu]$, by applying the subroutine Exploit-SI, given by Algorithm 4. This subroutine is extended from [12] based on our coded storage scenario. Note that in Algorithm 4 the function $\text{Reproduce}(j, Q^{(v)}_{j'}(u, \mathcal{U}; \tau - 1))$, $j' \in [1 : n] \setminus \{j\}$, simply reproduces all the queries in the auxiliary query set $Q^{(v)}_{j'}(u, \mathcal{U}; \tau - 1)$ with a different coordinate $j$.

**Algorithm 4: Exploit-SI**

> **Input** : $u$, $Q^{(v)}_1(u, \mathcal{U}; \tau - 1), \ldots, Q^{(v)}_n(u, \mathcal{U}; \tau - 1)$, $j$, $v$, and $\tau$
> **Output:** $\phi^{(v)}(u, \mathcal{U}; \tau - 1)$
> 1  $\phi^{(v)}(u, \mathcal{U}; \tau - 1) \leftarrow \emptyset$
> 2  **for** $i \in [1 : \kappa]$ **do**
> 3  $\quad$ **for** $j' \in [1 : n] \setminus \{j\}$ **do**
> 4  $\qquad$ **if** $u = a_{i,j'}$ **then**
> 5  $\qquad\quad \phi^{(v)}(u, \mathcal{U}; \tau - 1) \leftarrow \text{Reproduce}(j, Q^{(v)}_{j'}(u, \mathcal{U}; \tau - 1))$
> 6  $\qquad\quad$ break
> 7  $\qquad$ **end**
> 8  $\quad$ **end**
> 9  **end**

Next, we update the desired query sets $Q^{(v)}_j(a_{i,j}, \mathcal{D}; \tau)$ in Steps 25 to 31. First, the function $\text{Partition}\big(\tilde{Q}^{(v)}_j(\mathcal{U}; \tau - 1)\big)$ denotes a procedure that divides a set into $\kappa$ disjoint equally-sized subsets. This is viable since based on the subroutine Initial-Round and the following subroutine M-Sym, one can show that $|\tilde{Q}^{(v)}_j(\mathcal{U}; \tau - 1)| = \binom{\mu-1}{\tau-1} \kappa^{\mu-(\tau-1)} (\nu - \kappa)^{(\tau-1)-1} \cdot (\nu - \kappa)$ for each round

$\tau \in [2 : \mu]$, which is always divisible by $\kappa$. Secondly, we assign the new query set of desired symbols $Q^{(v)}_j(a_{i,j}, \mathcal{D}; \tau)$ for the current round by using an element-wise set addition $\text{SetAddition}(Q_1, Q_2)$. The element-wise set addition is defined as $\{q_{i_l} + q_{i'_l} : q_{i_l} \in Q_1, q_{i'_l} \in Q_2, l \in [1 : \rho]\}$ with $|Q_1| = |Q_2| = \rho$, where $\rho$ is an appropriate integer. In Steps 33 to 37, the subroutine M-Sym, given in Algorithm 5, is invoked to generate the undesired query sets $Q^{(v)}_j(a_{i,j}, \mathcal{U}; \tau)$ by utilizing message symmetry. This subroutine selects symbols of undesired messages to generate $\tau$-sums that enforce symmetry in the round queries. The procedure resembles the subroutine M-Sym proposed in [12]. In Algorithm 5, $\Pi_\tau$ denotes the set of all possible selections of $\tau$ distinct indices in $[1 : \mu]$ and $\text{Lexico}(\Pi_\tau)$ denotes the corresponding set of ordered selections (the indices $(v_1, \ldots, v_\tau)$ of a selection of $\Pi_\tau$ are ordered in natural lexicographical order). Further, the notation $U^{(v_x)}_{*,j}$ implies that the row index of the code symbol can be arbitrary. This is the case since only the linear function indices $(v_1, \ldots, v_\tau)$ are necessary to determine $i_z$, $\forall z \in [1 : \tau]$. As a result, symmetry over the linear functions is maintained.

**Algorithm 5: M-Sym**

> **Input** : $Q^{(v)}_j(u, \mathcal{D}; \tau)$, $j$, $v$, and $\tau$
> **Output:** $\phi^{(v)}(u, \mathcal{U}; \tau)$
> 1  $\phi^{(v)}(u, \mathcal{U}; \tau) \leftarrow \emptyset$
> 2  **for** $(v_1, \ldots, v_\tau) \in \text{Lexico}(\Pi_\tau)$, $v \notin \{v_1, \ldots, v_\tau\}$ **do**
> 3  $\quad \phi^{(v)}(u, \mathcal{U}; \tau) \leftarrow \phi^{(v)}(u, \mathcal{U}; \tau) \cup \{U^{(v_1)}_{i_1, j} + \ldots + U^{(v_\tau)}_{i_\tau, j}\}$
>    $\qquad$ such that $\forall z \in [1 : \tau]$,
>    $\qquad \exists U^{(v)}_{i_z, j} + \sum_{\substack{x \in [1:\tau] \\ x \neq z}} U^{(v_x)}_{*,j} \in Q^{(v)}_j(u, \mathcal{D}; \tau)$
> 4  **end**

*6) Query Set Assembly:* Finally, in Steps 39 to 48, we assemble each query set from disjoint query subsets obtained in all $\tau$ rounds. It can be shown that $Q^{(v)}_j(\mathcal{D}; \tau) \cup Q^{(v)}_j(\mathcal{U}; \tau)$ contains $\kappa^{\mu-(\tau-1)} (\nu - \kappa)^{\tau-1}$ $\tau$-sums for every $\tau$-sum type.

### B. Sign Assignment and Redundancy Elimination

After Algorithm 1, the user will know which row indices of the stored code symbols he/she is going to request. To reduce the total number of downloaded symbols, the linear dependency among the candidate linear functions is exploited. We adopt a similar sign assignment process over $\sigma^{(v)}_t \in \{+1, -1\}$ to each symbol in the query sets: $\sigma^{(v)}_t U^{(v)}_{t,j}$, based on the desired linear function index $v$ as introduced in [12, Sec. 4.2]. Here, the sign $\sigma^{(v)}_t$ is privately generated by the user with a uniform distribution over $\{-1, +1\}$. The intuition behind the sign assignment is to introduce a uniquely solvable equation system from the different $\tau$-sum types given the side information available from all other databases. By obtaining such a system of equations in each round, the user can determine some of the queries offline to decode the desired linear function and/or interference, thus reducing the download rate. Based on this insight we can state the following lemma.

*Lemma 1:* For all $v \in [1 : \mu]$, each database $j \in [1 : n]$, and based on the side information available from the databases, there are $\binom{\mu-r}{\tau}$ redundant $\tau$-sum types out of all possible types $\binom{\mu}{\tau}$ in each round $\tau \in [1 : \mu - r]$ of the query sets.

Since repetition codes and MDS codes are MDS-PIR capacity-achieving codes, Lemma 1 generalizes [12, Lem. 1] and [14, Lem. 1]. We can extend the proof of [12, Lem. 1] to our scheme based on the insight that the redundancy resulting from the linear dependencies between messages is also present for MDS-PIR capacity-achieving codes. We now make the final modification to our PLC query sets. We first directly apply the sign assignment $\sigma_t^{(v)}$, then remove all $\tau$-sums corresponding to redundant $\tau$-sum types from every round $\tau \in [1 : \mu]$ (see Lemma 1). Note that, in our case the redundancy is dependent on the rank of the functions matrix, $\mathrm{rank}(\mathsf{V}) = r \leq \min\{\mu, f\}$, thus generalizing the MDS-coded PLC case. Finally, we generate the queries $Q_{1:n}^{(v)}$.

### C. Example

To illustrate the key idea of the scheme we use the following example.

Consider four messages $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$, $\mathbf{W}^{(3)}$, and $\mathbf{W}^{(4)}$ that are stored in a DSS using a $[4, 2]$ MDS-PIR capacity-achieving code $\mathscr{C}$ with generator matrix

$$\mathsf{G}^{\mathscr{C}} = \begin{pmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{pmatrix} \text{ for which } \Lambda_{1,2} = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}$$

is a PIR achievable rate matrix. According to Definition 5 we obtain the PIR interference matrices $\mathsf{A}_{1 \times 4} = \begin{pmatrix} 1 & 2 & 1 & 2 \end{pmatrix}$ and $\mathsf{B}_{1 \times 4} = \begin{pmatrix} 2 & 1 & 2 & 1 \end{pmatrix}$. Suppose that the user wishes to obtain a linear function $\mathbf{X}^{(v)}$ from a set of $\mu = 3$ candidate linear functions, i.e., $v \in [1 : 3]$, defined by

$$\mathsf{V}_{\mu \times f} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 2 & 1 & 0 & 1 \end{pmatrix}.$$

We simplify notation by letting $x_{t,j} = U_{t,j}^{(1)}$, $y_{t,j} = U_{t,j}^{(2)}$, and $z_{t,j} = U_{t,j}^{(3)}$ for all $t \in [1 : \beta]$, $j \in [1 : n]$, where $\beta = \nu^\mu = 8$. Let the desired linear function index be $v = 1$. For this example, the construction of the query sets is briefly presented in the following steps.

***Initialization (Round $\tau = 1$):*** Algorithm 1 starts with $\tau = 1$ to generate auxiliary query sets for each database holding $\kappa^{\mu-1} = 1$ distinct instances of $x_{t,j}$. By message symmetry this also applies to $y_{t,j}$ and $z_{t,j}$. The auxiliary query sets for the first round are given in Table I(a). Note that the queries for $z_{t,j}$ can be generated offline by the user and thus are later removed from the query sets.

***Following Rounds ($\tau > 1$):*** As can be seen from Table I(b) and (c), using the interference matrices $\mathsf{A}_{1 \times 4}$ and $\mathsf{B}_{1 \times 4}$, the algorithm further generates the auxiliary query sets $Q_j^{(1)}(a_{1,j}, \mathcal{D}; \tau)$ that contain new symbols to be queried jointly with symbols of the desired linear function during the exploitation of side information for the $j$-th database, $j \in [1 : n]$. After generating the desired auxiliary query

sets $Q_j^{(1)}(a_{1,j}, \mathcal{D}; \tau)$, the undesired auxiliary query sets $Q_j^{(1)}(a_{1,j}, \mathcal{U}; \tau)$ are generated by enforcing message symmetry. Finally, we apply the sign assignment scheme to the query sets and remove redundant queries from each database. The resulting queries are shown in Table II.

The PLC rate of the achievable scheme is equal to $\frac{\nu^\mu \cdot k}{\mathsf{D}} = \frac{8 \cdot 2}{6 \cdot 4} = \frac{2}{3}$, which is equal to the PLC capacity in (2) with $\mathrm{rank}(\mathsf{V}) = 2$.

TABLE I
AUXILIARY QUERY SETS FOR EACH ROUND. THE MAGENTA DASHED ARROWS AND THE CYAN ARROWS INDICATE THAT THE EXPLOIT-SI ALGORITHM AND THE M-SYM ALGORITHM ARE USED, RESPECTIVELY.

| $j$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $Q_j^{(1)}(a_{1,j}, \mathcal{D}; 1)$ | $x_{(1-1)\cdot 1+1,1}$ | $x_{(2-1)\cdot 1+1,2}$ | $x_{(1-1)\cdot 1+1,3}$ | $x_{(2-1)\cdot 1+1,4}$ |
| $Q_j^{(1)}(a_{1,j}, \mathcal{U}; 1)$ | $y_{(1-1)\cdot 1+1,1}$ | $y_{(2-1)\cdot 1+1,2}$ | $y_{(1-1)\cdot 1+1,3}$ | $y_{(2-1)\cdot 1+1,4}$ |
| | $z_{(1-1)\cdot 1+1,1}$ | $z_{(2-1)\cdot 1+1,2}$ | $z_{(1-1)\cdot 1+1,3}$ | $z_{(2-1)\cdot 1+1,4}$ |

(a)

| $j$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $Q_j^{(1)}(a_{1,j}, \mathcal{D}; 2)$ | $x_{1\cdot 2+1,1} + y_{2,1}$ | $x_{1\cdot 2+2,2} + y_{1,2}$ | $x_{1\cdot 2+1,3} + y_{2,3}$ | $x_{1\cdot 2+2,4} + y_{1,4}$ |
| | $x_{2\cdot 2+1,1} + z_{2,1}$ | $x_{2\cdot 2+2,2} + z_{1,2}$ | $x_{2\cdot 2+1,3} + z_{2,3}$ | $x_{2\cdot 2+2,4} + z_{1,4}$ |
| $Q_j^{(1)}(a_{1,j}, \mathcal{U}; 2)$ | $y_{4+1,1} + z_{2+1,1}$ | $y_{4+2,2} + z_{2+2,2}$ | $y_{4+1,3} + z_{2+1,3}$ | $y_{4+2,4} + z_{2+2,4}$ |

(b)

| $j$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $Q_j^{(1)}(a_{1,j}, \mathcal{D}; 3)$ | $x_{3\cdot 2+1,1} + y_{6,1} + z_{4,1}$ | $x_{3\cdot 2+1,2} + y_{5,2} + z_{3,2}$ | $x_{3\cdot 2+1,3} + y_{6,3} + z_{4,3}$ | $x_{3\cdot 2+1,4} + y_{5,4} + z_{3,4}$ |

(c)

TABLE II
QUERY SETS FOR A CODED DSS THAT USES THE $[4, 2]$ CODE OF EXAMPLE IV-C TO STORE $f = 4$ MESSAGES.

| $j$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $Q_j^{(1)}(\mathcal{D}; 1)$ | $x_{1,1}$ | $x_{2,2}$ | $x_{1,3}$ | $x_{2,4}$ |
| $Q_j^{(1)}(\mathcal{U}; 1)$ | $y_{1,1}$ | $y_{2,2}$ | $y_{1,3}$ | $y_{2,4}$ |
| $Q_j^{(1)}(\mathcal{D}; 2)$ | $x_{3,1} + y_{2,1}$ | $x_{4,2} + y_{1,2}$ | $x_{3,3} + y_{2,3}$ | $x_{4,4} + y_{1,4}$ |
| | $x_{5,1} + z_{2,1}$ | $x_{6,2} + z_{1,2}$ | $x_{5,3} + z_{2,3}$ | $x_{6,4} + z_{1,4}$ |
| $Q_j^{(1)}(\mathcal{U}; 2)$ | $y_{5,1} + z_{3,1}$ | $y_{6,2} + z_{4,2}$ | $y_{5,3} + z_{3,3}$ | $y_{6,4} + z_{4,4}$ |
| $Q_j^{(1)}(\mathcal{D}; 3)$ | $x_{7,1} + y_{6,1} + z_{4,1}$ | $x_{8,2} + y_{5,2} + z_{3,2}$ | $x_{7,3} + y_{6,3} + z_{4,3}$ | $x_{8,4} + y_{5,4} + z_{3,4}$ |

### D. Achievable Rate

The resulting achievable rate of Algorithm 1 after removing redundant $\tau$-sums according to Lemma 1 is given as

$$\mathsf{R} \overset{(a)}{=} \frac{k\nu^\mu}{n \sum_{\tau=1}^{\mu} \left( \binom{\mu}{\tau} - \binom{\mu-r}{\tau} \right) \kappa^{\mu-(\tau-1)} (\nu - \kappa)^{\tau-1}}$$

$$\overset{(b)}{=} \frac{\kappa\nu^\mu}{\nu \sum_{\tau=1}^{\mu} \left( \binom{\mu}{\tau} - \binom{\mu-r}{\tau} \right) \kappa^{\mu-(\tau-1)} (\nu - \kappa)^{\tau-1}}$$

$$= \frac{\nu^\mu \left( \frac{\nu-\kappa}{\nu} \right)}{\sum_{\tau=1}^{\mu} \left( \binom{\mu}{\tau} - \binom{\mu-r}{\tau} \right) \kappa^{\mu-\tau} (\nu - \kappa)^{\tau}}$$

$$\vdots$$

$$\overset{(c)}{=} \frac{\nu^\mu \left( 1 - \frac{\kappa}{\nu} \right)}{\nu^\mu - \kappa^r \nu^{\mu-r}}$$

$$= \left( 1 - \frac{\kappa}{\nu} \right) \left[ 1 - \left( \frac{\kappa}{\nu} \right)^r \right]^{-1},$$

where we define $\binom{m}{n} \triangleq 0$ if $m < n$; (a) follows from the PLC rate in Definition 2, the fact that $Q_j^{(v)}(\mathcal{D}; \tau) \cup Q_j^{(v)}(\mathcal{U}; \tau)$

contains $\kappa^{\mu-(\tau-1)}(\nu-\kappa)^{\tau-1}$ $\tau$-sums for every $\tau$-sum type, and Lemma 1; (b) follows from Definition 4; and (c) follows by adapting similar steps as in the proof given in [14].

## V. Converse Proof

In [18], [19], the PIR capacity of MDS-PIR capacity-achieving codes is shown to be equal to the MDS-PIR capacity. In this section, we adapt the converse proof of [19, Thm. 4] to the scenario of the linearly-coded PLC problem. We will show that the PLC capacity of MDS-PIR capacity-achieving codes is equal to (2). Before we proceed with the converse proof, we give some general results that are useful for the proof.

1) From the condition of privacy,

$$\mathsf{H}\big(A_j^{(v)} \,\big|\, \mathbf{X}^{(v)}, \mathcal{Q}\big) = \mathsf{H}\big(A_j^{(v')} \,\big|\, \mathbf{X}^{(v)}, \mathcal{Q}\big), \quad (3)$$

where $v \neq v'$, $v, v' \in [1:\mu]$, and $\mathcal{Q} \triangleq \big\{Q_j^{(v)} \colon v \in [1:\mu], j \in [1:n]\big\}$ denotes the set of all possible queries made by the user. Although this seems to be intuitively true, a proof of this property is still required and can be found in [21, Lem. 3].

2) Consider a PLC protocol for a coded DSS that uses an $[n,k]$ code $\mathscr{C}$ to store $f$ messages. For any subset of linear combinations $\mathcal{V} \subseteq [1:\mu]$ and for any information set $\mathcal{I}$ of $\mathscr{C}$, we have

$$\mathsf{H}\big(A_{\mathcal{I}}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) = \sum_{j \in \mathcal{I}} \mathsf{H}\big(A_j^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big). \quad (4)$$

The proof uses the linear independence of the columns of a generator matrix of $\mathscr{C}$ corresponding to an information set, and can be seen as a simple extension of [8, Lem. 2] or [21, Lem. 4]. This argument applies to the case of PLC due to the fact that $A_{\mathcal{I}}^{(v)}$ is still a deterministic function of independent random variables $\{\boldsymbol{C}_j \colon j \in \mathcal{I}\}$ and $\mathcal{Q}$, where $\boldsymbol{C}_j \triangleq \big(C_{1,j}^{(1)}, \ldots, C_{\beta,j}^{(1)}, C_{1,j}^{(2)}, \ldots, C_{\beta,j}^{(f)}\big)^\mathsf{T}$ denotes the $f$ coded chunks that are stored in the $j$-th database.

Next, we state Shearer's Lemma, which represents a very useful entropy method for combinatorial problems.

*Lemma 2 (Shearer's Lemma [22]):* Let $\mathscr{S}$ be a collection of subsets of $[1:n]$, with each $j \in [1:n]$ included in at least $\kappa$ members of $\mathscr{S}$. For random variables $Z_1, \ldots, Z_n$, we have

$$\sum_{\mathcal{S} \in \mathscr{S}} \mathsf{H}(Z_{\mathcal{S}}) \geq \kappa \,\mathsf{H}(Z_1, \ldots, Z_n).$$

For our converse proof for the coded PLC problem, we also need the following lemma, whose proof has been omitted due to lack of space.

*Lemma 3:* Consider the linear mapping $\mathsf{V} = (v_{i,j})$ defined in (1) with $\mathrm{rank}(\mathsf{V}) = r$ where $v_{i_1,j_1}, \ldots, v_{i_r,j_r}$ are the entries corresponding to the pivot elements of $\mathsf{V}$. It follows that $\big(\mathbf{X}^{(i_1)}, \ldots, \mathbf{X}^{(i_h)}\big)$ and $\big(\mathbf{W}^{(j_1)}, \ldots, \mathbf{W}^{(j_h)}\big)$ are identically distributed, for some $h \in [1:r]$. In other words, $\mathsf{H}\big(\mathbf{X}^{(i_1)}, \ldots, \mathbf{X}^{(i_h)}\big) = h\mathsf{L}$, $h \in [1:r]$.

Now, we are ready for the converse proof. By [17, Lem. 2], since the code $\mathscr{C}$ is MDS-PIR capacity-achieving,

there exist $\nu$ information sets $\mathcal{I}_1, \ldots, \mathcal{I}_\nu$ such that each coordinate $j \in [1:n]$ is included in exactly $\kappa$ members of $\mathscr{I} = \{\mathcal{I}_1, \ldots, \mathcal{I}_\nu\}$ with $\frac{\kappa}{\nu} = \frac{k}{n}$.

Applying the chain rule of entropy we have

$$\mathsf{H}\big(A_{[1:n]}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) \geq \mathsf{H}\big(A_{\mathcal{I}_i}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big), \quad \forall\, i \in [1:\nu].$$

Let $v \in \mathcal{V}$ and $v' \in \mathcal{V}^{\mathrm{c}} \triangleq [1:\mu] \setminus \mathcal{V}$. Following similar steps as in the proof given in [8], [21], we get

$$\nu\,\mathsf{H}\big(A_{[1:n]}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)$$

$$\geq \sum_{i=1}^{\nu} \mathsf{H}\big(A_{\mathcal{I}_i}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)$$

$$= \sum_{i=1}^{\nu} \left( \sum_{j \in \mathcal{I}_i} \mathsf{H}\big(A_j^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) \right) \quad (5)$$

$$= \sum_{i=1}^{\nu} \left( \sum_{j \in \mathcal{I}_i} \mathsf{H}\big(A_j^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) \right) \quad (6)$$

$$= \sum_{i=1}^{\nu} \mathsf{H}\big(A_{\mathcal{I}_i}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) \quad (7)$$

$$\geq \kappa\,\mathsf{H}\big(A_{[1:n]}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big) \quad (8)$$

$$= \kappa\Big[\mathsf{H}\big(A_{[1:n]}^{(v')}, \mathbf{X}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)$$

$$\qquad - \mathsf{H}\big(\mathbf{X}^{(v')} \,\big|\, A_{[1:n]}^{(v')}, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)\Big]$$

$$= \kappa\Big[\mathsf{H}\big(\mathbf{X}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)$$

$$\qquad + \mathsf{H}\big(A_{[1:n]}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathbf{X}^{(v')}, \mathcal{Q}\big) - 0\Big] \quad (9)$$

$$= \kappa\Big[\mathsf{H}\big(\mathbf{X}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}\big) + \mathsf{H}\big(A_{[1:n]}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathbf{X}^{(v')}, \mathcal{Q}\big)\Big], \quad (10)$$

where (5) and (7) follow from (4); (6) is because of (3); (8) is due to Shearer's Lemma; (9) is from the fact that the $v'$-th linear combination $\mathbf{X}^{(v')}$ is determined by the answers $A_{[1:n]}^{(v')}$ and all possible queries $\mathcal{Q}$; and finally, (10) follows from the independence between all possible queries and the messages. Therefore, we can conclude that

$$\mathsf{H}\big(A_{[1:n]}^{(v)} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathcal{Q}\big)$$

$$\geq \frac{\kappa}{\nu}\,\mathsf{H}\big(\mathbf{X}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}\big) + \frac{\kappa}{\nu}\,\mathsf{H}\big(A_{[1:n]}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathbf{X}^{(v')}, \mathcal{Q}\big)$$

$$= \frac{k}{n}\,\mathsf{H}\big(\mathbf{X}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}\big) + \frac{k}{n}\,\mathsf{H}\big(A_{[1:n]}^{(v')} \,\big|\, \mathbf{X}^{\mathcal{V}}, \mathbf{X}^{(v')}, \mathcal{Q}\big), \quad (11)$$

where we have used Definition 4 to obtain (11).

Since there are in total $\mu$ linear combinations and $\mathcal{L} \triangleq \{\ell_1, \ldots, \ell_r\} \subseteq [1:\mu]$ is the set of row indices corresponding to the selected basis for the row space of $\mathsf{V}$, we can recursively use (11) $r-1$ times to obtain

$$\mathsf{H}\big(A_{[1:n]}^{(\ell_1)} \,\big|\, \mathbf{X}^{(\ell_1)}, \mathcal{Q}\big)$$

$$\geq \sum_{v=1}^{r-1} \Big(\frac{k}{n}\Big)^v \mathsf{H}\big(\mathbf{X}^{(\ell_{v+1})} \,\big|\, \mathbf{X}^{\{\ell_1, \ldots, \ell_v\}}\big)$$

$$\qquad + \Big(\frac{k}{n}\Big)^{r-1} \mathsf{H}\big(A_{[1:n]}^{(\ell_r)} \,\big|\, \mathbf{X}^{\{\ell_1, \ldots, \ell_r\}}, \mathcal{Q}\big)$$

$$\geq \sum_{v=1}^{r-1} \left(\frac{k}{n}\right)^v \mathsf{H}\big(\mathbf{X}^{(\ell_{v+1})} \,\big|\, \mathbf{X}^{\{\ell_1,\ldots,\ell_v\}}\big) \tag{12}$$

$$= \sum_{v=1}^{r-1} \left(\frac{k}{n}\right)^v \mathsf{L}, \tag{13}$$

where (12) follows from the nonnegativity of entropy. (13) holds since it follows from Lemma 3 that $\mathsf{H}\big(\mathbf{X}^{(\ell_{v+1})} \,\big|\, \mathbf{X}^{\{\ell_1,\ldots,\ell_v\}}\big) = \mathsf{H}\big(\mathbf{X}^{(\ell_{v+1})}\big) = \mathsf{L}$. Here, we also remark that the recursive steps follow the same principle of the general converse for dependent PIR (DPIR) from [15, Thm. 1]. In [15], the authors claim that the general converse for the DPIR problem strongly depends on the chosen permutation of the indices of the candidate functions. However, the best index permutation for the candidate linear functions for the PLC problem results from finding a basis for $\mathsf{V}$. Now,

$$\mathsf{L} = \mathsf{H}\left(\mathbf{X}^{(\ell_1)}\right)$$

$$= \mathsf{H}\big(\mathbf{X}^{(\ell_1)} \,\big|\, \mathcal{Q}\big) - \underbrace{\mathsf{H}\big(\mathbf{X}^{(\ell_1)} \,\big|\, A_{[1:n]}^{(\ell_1)}, \mathcal{Q}\big)}_{=0} \tag{14}$$

$$= \mathsf{I}\big(\mathbf{X}^{(\ell_1)}; A_{[1:n]}^{(\ell_1)} \,\big|\, \mathcal{Q}\big)$$

$$= \mathsf{H}\left(A_{[1:n]}^{(\ell_1)} \,\Big|\, \mathcal{Q}\right) - \mathsf{H}\left(A_{[1:n]}^{(\ell_1)} \,\Big|\, \mathbf{X}^{(\ell_1)}, \mathcal{Q}\right)$$

$$\leq \mathsf{H}\left(A_{[1:n]}^{(\ell_1)} \,\Big|\, \mathcal{Q}\right) - \sum_{v=1}^{r-1} \left(\frac{k}{n}\right)^v \mathsf{L}, \tag{15}$$

where (14) follows since any message is independent of the queries $\mathcal{Q}$, and knowing the answers $A_{[1:n]}^{(\ell_1)}$ and the queries $\mathcal{Q}$, one can determine $\mathbf{X}^{(\ell_1)}$; (15) holds because of (13).

Finally, the converse proof is completed by showing that

$$\mathsf{R} = \frac{\mathsf{L}}{\sum_{j=1}^n \mathsf{H}\left(A_j^{(\ell_1)}\right)}$$

$$\leq \frac{\mathsf{L}}{\mathsf{H}\left(A_{[1:n]}^{(\ell_1)}\right)} \tag{16}$$

$$\leq \frac{\mathsf{L}}{\mathsf{H}\left(A_{[1:n]}^{(\ell_1)} \,\Big|\, \mathcal{Q}\right)} \tag{17}$$

$$\leq \frac{1}{1 + \sum_{v=1}^{r-1} \left(\frac{k}{n}\right)^v} = \mathsf{C}_{\mathrm{PLC}}, \tag{18}$$

where (16) holds because of the chain rule of entropy, (17) is due to the fact that conditioning reduces entropy, and we apply (15) to obtain (18).

## VI. CONCLUSION

We have provided the capacity of PLC for DSSs, where data is encoded and stored using an arbitrary linear code from the class of MDS-PIR capacity-achieving codes. Interestingly, the capacity is equal to the corresponding PIR capacity. Thus, privately retrieving arbitrary linear combinations of the stored messages does not incur any overhead in rate compared to retrieving a single message from the databases.

## REFERENCES

[1] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, "Private information retrieval," *J. ACM*, vol. 45, no. 6, pp. 965–981, Nov. 1998.

[2] S. Yekhanin, "Private information retrieval," *Commun. ACM*, vol. 53, no. 4, pp. 68–73, Apr. 2010.

[3] N. B. Shah, K. V. Rashmi, and K. Ramchandran, "One extra bit of download ensures perfectly private information retrieval," in *Proc. IEEE Int. Symp. Inf. Theory*, Honolulu, HI, USA, Jun. 29 – Jul. 4, 2014, pp. 856–860.

[4] T. H. Chan, S.-W. Ho, and H. Yamamoto, "Private information retrieval for coded storage," in *Proc. IEEE Int. Symp. Inf. Theory*, Hong Kong, China, Jun. 14–19, 2015, pp. 2842–2846.

[5] R. Tajeddine, O. W. Gnilke, and S. El Rouayheb, "Private information retrieval from MDS coded data in distributed storage systems," 2018, to app. in *IEEE Trans. Inf. Theory*.

[6] H. Sun and S. A. Jafar, "The capacity of private information retrieval," *IEEE Trans. Inf. Theory*, vol. 63, no. 7, pp. 4075–4088, Jul. 2017.

[7] ——, "The capacity of robust private information retrieval with colluding databases," *IEEE Trans. Inf. Theory*, vol. 64, no. 4, pp. 2361–2370, Apr. 2018.

[8] K. Banawan and S. Ulukus, "The capacity of private information retrieval from coded databases," *IEEE Trans. Inf. Theory*, vol. 64, no. 3, pp. 1945–1956, Mar. 2018.

[9] R. Freij-Hollanti, O. W. Gnilke, C. Hollanti, and D. A. Karpuk, "Private information retrieval from coded databases with colluding servers," *SIAM J. Appl. Algebra Geom.*, vol. 1, no. 1, pp. 647–664, Nov. 2017.

[10] S. Kumar, E. Rosnes, and A. Graell i Amat, "Private information retrieval in distributed storage systems using an arbitrary linear code," in *Proc. IEEE Int. Symp. Inf. Theory*, Aachen, Germany, Jun. 25–30, 2017, pp. 1421–1425.

[11] M. Mirmohseni and M. A. Maddah-Ali, "Private function retrieval," Nov. 2017, arXiv:1711.04677v2 [cs.IT]. [Online]. Available: https://arxiv.org/abs/1711.04677

[12] H. Sun and S. A. Jafar, "The capacity of private computation," Oct. 2017, arXiv:1710.11098v3 [cs.IT]. [Online]. Available: https://arxiv.org/abs/1710.11098

[13] D. Karpuk, "Private computation of systematically encoded data with colluding servers," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, Jun. 17–22, 2018, pp. 2112–2116.

[14] S. A. Obead and J. Kliewer, "Achievable rate of private function retrieval from MDS coded databases," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, Jun. 17–22, 2018, pp. 2117–2121.

[15] Z. Chen, Z. Wang, and S. Jafar, "The asymptotic capacity of private search," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, Jun. 17–22, 2018, pp. 2122–2126.

[16] S. Kumar, H.-Y. Lin, E. Rosnes, and A. Graell i Amat, "Achieving maximum distance separable private information retrieval capacity with linear codes," Dec. 2017, arXiv:1712.03898v3 [cs.IT]. [Online]. Available: https://arxiv.org/abs/1712.03898

[17] H.-Y. Lin, S. Kumar, E. Rosnes, and A. Graell i Amat, "An MDS-PIR capacity-achieving protocol for distributed storage using non-MDS linear codes," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, Jun. 17–22, 2018, pp. 966–970.

[18] ——, "Asymmetry helps: Improved private information retrieval protocols for distributed storage," to be pres. at *IEEE Inf. Theory Workshop*, Guangzhou, China, Nov. 25–29, 2018. [Online]. Available: https://arxiv.org/abs/1806.01342

[19] ——, "On the fundamental limit of private information retrieval for coded distributed storage," Aug. 2018, arXiv:1808.09018v1 [cs.IT]. [Online]. Available: https://arxiv.org/abs/1808.09018

[20] R. Freij-Hollanti, O. W. Gnilke, C. Hollanti, A.-L. Horlemann-Trautmann, D. Karpuk, and I. Kubjas, "*t*-private information retrieval schemes using transitive codes," Dec. 2017, arXiv:1712.02850v1 [cs.IT]. [Online]. Available: https://arxiv.org/abs/1712.02850

[21] J. Xu and Z. Zhang, "On sub-packetization of capacity-achieving PIR schemes for MDS coded databases," Dec. 2017, arXiv:1712.02466v2 [cs.IT]. [Online]. Available: http://arxiv.org/abs/1712.02466

[22] J. Radhakrishnan, "Entropy and counting," in *Computational mathematics, modelling and algorithms*, J. C. Misra, Ed. Narosa Publishing House, 2003, pp. 146–168.