# Secure Private Information Retrieval from Colluding Databases with Eavesdroppers

Qiwen Wang, and Mikael Skoglund

School of Electrical Engineering, KTH Royal Institute of Technology

Email: {qiwenw, skoglund}@kth.se

## Abstract

The problem of *private information retrieval (PIR)* is to retrieve one message out of $K$ messages replicated at $N$ databases, without revealing the identity of the desired message to the databases. We consider the problem of PIR with colluding servers and eavesdroppers, named *T-EPIR*. Specifically, any $T$ out of $N$ databases may collude, that is, they may communicate their interactions with the user to guess the identity of the requested message. An eavesdropper is curious to know the database and can tap in on the incoming and outgoing transmissions of any $E$ databases. The databases share some common randomness unknown to the eavesdropper and the user, and use the common randomness to generate the answers, such that the eavesdropper can learn no information about the $K$ messages. Define $R^*$ as the optimal ratio of the number of the desired message information bits to the number of total downloaded bits, and $\rho^*$ to be the optimal ratio of the information bits of the shared common randomness to the information bits of the desired file. In our previous work [1], we found that when $E \geq T$, the optimal ratio that can be achieved (hence is the capacity) equals $1 - \frac{E}{N}$. In this work, we focus on the case when $E \leq T$. We derive an outer bound (converse bound) that $R^* \leq \left(1 - \frac{T}{N}\right) \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K}$. We also obtain a lower bound (converse bound) of $\rho^* \geq \frac{\frac{E}{N}\left(1 - \left(\frac{T}{N}\right)^K\right)}{\left(1 - \frac{T}{N}\right)\left(1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}\right)}$. For the achievability, we propose a scheme which achieves the rate (inner bound) $R = \frac{1 - \frac{T}{N}}{1 - \left(\frac{T}{N}\right)^K} - \frac{E}{KN}$. The amount of shared common randomness used in the achievable scheme is $\frac{\frac{E}{N}\left(1 - \left(\frac{T}{N}\right)^K\right)}{1 - \frac{T}{N} - \frac{E}{KN}\left(1 - \left(\frac{T}{N}\right)^K\right)}$ times the file size. The gap between the derived inner and outer bounds vanishes as the number of messages $K$ tends to infinity.

## I. Introduction

In the situation where a user wants to retrieve a file (message) from a remotely stored database, the nature of the data might be privacy-sensitive, for example medical records, stock prices *etc.*, such that the user does not want to reveal the identity of the data retrieved. This is known as

the problem of private information retrieval (PIR). In some cases, the privacy of the database needs also to be preserved. For example, if a user wants to retrieve his/her medical data from a database, it is hoped that the user obtains no information about other users' medical records. This is known as the problem of symmetric private information retrieval (SPIR).

The problem of PIR and SPIR was firstly studied in the computer science literature. In [2], [3], it is shown that if the messages are stored at a single database, the only possible scheme for the user is to download all the messages to guarantee information-theoretic privacy, which is inefficient in practice. It is further shown that the communication cost can be reduced in sublinear scale by replicating the database at multiple non-colluding servers [3]. To further protect the privacy of the database such that the user obtains no more information regarding the other messages besides the requested message, the problem of SPIR is introduced [4]. In [2]–[4], the collection of messages stored at each database is modeled as a bit string, and the user wishes to retrieve a single bit. In these works, the communication cost is measured as the sum of the transmission at the querying phase from user to servers and at the downloading phase from servers to user.

When the message size is significantly large and the target is to minimize the communication cost of only the downloading phase, the metric of the downloading cost is defined as the number of bits downloaded per bit of the retrieved message, and the reciprocal of which is named the *PIR capacity*. A series of recent works derive information-theoretic limits of various versions of the PIR problem [5]–[11] *etc*. The leading work in the area is by Sun and Jafar [5], where the authors find the capacity of the PIR problem with replicated databases. In subsequent works by Sun and Jafar [6], [7], the PIR capacity with duplicated databases and colluding servers, and the SPIR capacity with duplicated (non-colluding) databases are derived. In [8]–[10], Banawan and Ulukus find the capacity of the PIR problem with coded databases, multi-message PIR with replicated databases, and the PIR problem with colluding and Byzantine databases. In our previous works [1], [11], [12], we derive the capacity of the SPIR problem with coded databases, linear SPIR with colluding and coded databases, and the SPIR problem with Byzantine adversaries and eavesdroppers.

Another series of works focus more on the coding structure of the storage system, and study schemes and information limits for various PIR problems with coded databases [13]–[17]. In [13], PIR is achieved by downloading one extra bit other than the desired file, given that the number

of storage nodes grows with file size, which can be impractical in some storage systems. In [14], storage overhead can be reduced by increasing the number of storage nodes. In [15], tradeoff between storage cost and downloading cost is analyzed. Subsequently in [16], explicit schemes which match the tradeoff in [15] are presented. It is worth noting that in [8], the capacity of PIR for coded database is settled, which improves the results in [15], [16]. Recently in [17], the authors present a framework for PIR from coded databases with colluding servers.

In our previous work [1], we studied the problem of SPIR from replicated databases with colluding databases and eavesdroppers, named *T-ESPIR*. Briefly speaking, a user wants to retrieve one file out of $K$ files that are replicatively stored at $N$ databases. Any $T$ out of the $N$ servers may collude, that is, they may share their communication with the user to infer the identity of the requested file. A passive eavesdropper is curious to know the database and can tap in on the incoming and outgoing transmissions of any $E$ servers. In the problem of T-ESPIR, it is required that the user learns no information about the database other than the requested file. In [1], we show that the information-theoretical capacity of the T-ESPIR problem is $1 - \frac{\max(T,E)}{N}$, if the databases share common randomness with amount at least $\frac{\max(T,E)}{N - \max(T,E)}$ times the file size. In Section VI.B in [1], we discussed that if database-privacy is not required, *i.e.* the user can learn information about the other files, and when $E \geq T$, the capacity of the T-EPIR problem is $1 - \frac{E}{N}$.

In this work, we continue the study of the T-EPIR problem when $E \leq T$. We derive an outer bound (converse bound) that $R^* \leq \left(1 - \frac{T}{N}\right) \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K}$. We also obtain a lower bound (converse bound) of $\rho^* \geq \frac{\frac{E}{N}\left(1 - \left(\frac{T}{N}\right)^K\right)}{\left(1 - \frac{T}{N}\right)\left(1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}\right)}$. For the achievability, we propose a scheme which achieves the rate (inner bound) $R = \frac{1 - \frac{T}{N}}{1 - \left(\frac{T}{N}\right)^K} - \frac{E}{KN}$. The amount of shared common randomness used in the achievable scheme is $\frac{\frac{E}{N}\left(1 - \left(\frac{T}{N}\right)^K\right)}{1 - \frac{T}{N} - \frac{E}{KN}\left(1 - \left(\frac{T}{N}\right)^K\right)}$ times the file size. The capacity of T-ESPIR when $E < T$ remains an open problem. In Section III, we discuss four special cases in which the capacity is known or can be easily derived, and reveal that our outer bound is tight for the four special cases. On the other hand, the inner bound is tight for three cases but one, namely, when $E = T$, the derived inner bound does not match with the capacity at this point. For illustration, we plot the results in Figure 1 and Figure 2 for some chosen parameters. It can be observed from the figures that the gap between inner and outer bounds decays and vanishes as $K$ tends to infinity.

## II. MODEL

### A. Notation

Let $[m : n]$ denote the set $\{m, m+1, \ldots, n\}$ for $m \leq n$. To simplify the notation, denote the set of random variables $\{X_m, X_{m+1}, \ldots, X_n\}$ by $X_{[m:n]}$. For an index set $\mathcal{I} = \{i_1, i_2, \ldots, i_n\}$, denote the set of variables with the index set $\{X_i : i \in \mathcal{I}\}$ by $X_{\mathcal{I}}$. For a matrix $\mathbf{S}$, let $\mathbf{S}[:, \mathcal{I}]$ denote the submatrix of $\mathbf{S}$ comprised of the columns corresponding to the index set $\mathcal{I}$. The transpose of matrix $\mathbf{G}$ is denoted by $\mathbf{G}^{\mathsf{T}}$. Let $\sim$ denote the statistical equivalence between random variables, that is, if $X \sim Y$, then $X$ and $Y$ are identically distributed.

### B. Problem Description

**Replicated databases:** A collection of $K$ independent messages (files), denoted by $W_1, \ldots, W_K$, are replicatively stored at $N$ databases (nodes). Each message consists $L$ information bits. Therefore, for any $k \in [1 : K]$,

$$H(W_k) = L \quad ; \quad H(W_1, \ldots, W_K) = KL.$$

**User queries:** A user wants to retrieve a message $W_\kappa$ with index $\kappa$ from the database, where the desired message index $\kappa$ follows some prior distribution among $[1 : K]$. Let $\mathcal{U}$ denote a random variable privately generated by the user, which represents the randomness of the query scheme followed by the user. The random variable $\mathcal{U}$ is generated independently of the messages and the desired file index. Let the realization of the file index $\kappa$ be $k$, based on the realization of the desired file index $k$ and the realization of $\mathcal{U}$, the user generates and sends queries to all nodes, where the query received by node-$n$ is denoted by $Q_n^{[k]}$. Let $\mathcal{Q} = [Q_n^{[k]}]_{n \in [1:N], k \in [1:K]}$ denote the complete query scheme, namely, the collection of all queries under all cases of desired message index. We have that $H(\mathcal{Q}|\mathcal{U}) = 0$.

**Common randomness:** Let random variable $S$ denote the common randomness shared by all databases, the realization of which is known to all the databases but unavailable to the user and the eavesdropper. The common randomness is utilized to protect the system-privacy (2) below, that is, to prevent the eavesdropper from learning the messages.

**Database answers:** The databases generate answers according to the agreed scheme with the user based on the received query $Q_n^{[k]}$, the stored messages $W_{[1:K]}$, and the common randomness

$S$. The answer generated and sent to the user by node $n$ is denoted by $A_n^{[k]}$.

**Eavesdropper:** A *passive eavesdropper* can tap in on the incoming and outgoing transmissions of $E$ nodes in the system. The eavesdropper is "nice but curious," in the sense that the goal of the eavesdropper is to obtain some information about the database, without corrupting any transmission. The user has no knowledge of the identity of the nodes tapped on by the eavesdropper.

**T-EPIR:** Based on the received answers $A_{[1:N]}^{[k]}$ and the query scheme $\mathcal{Q}$, the user shall be able to decode the requested message $W_k$ with zero error. Any set of $T$ databases may collude to guess the requested message index, by communicating their interactions with the user. Two privacy constraints must be satisfied:

- *User-privacy:* any $T$ colluding databases shall not be able to obtain any information regarding the identity of the requested message, *i.e.,*

$$I(\kappa; Q_{\mathcal{T}}^{[\kappa]}, A_{\mathcal{T}}^{[\kappa]}, W_{[1:K]}, S) = 0, \forall \mathcal{T} \subset [1:N], |\mathcal{T}| = T. \tag{1}$$

- *System-privacy:* For any set of databases $\mathcal{E}$ with size at most $E$, and for any $k \in [1:K]$:

$$I(W_{[1:K]}; Q_{\mathcal{E}}^{[k]}, A_{\mathcal{E}}^{[k]}) = 0. \tag{2}$$

**Definition 1.** *The rate of a T-EPIR scheme is the number of information bits of the requested file retrieved per downloaded answer bit. By symmetry among all files, for any $k \in [1:K]$,*

$$R_{\text{T-EPIR}} \triangleq \frac{H(W_k)}{\sum_{n=1}^{N} H(A_n^{[k]})}.$$

*The optimal rate of T-EPIR schemes is denoted by $R_{\text{T-EPIR}}^*$. The capacity $C_{\text{T-EPIR}}$ is the supremum of $R_{\text{T-EPIR}}$ over all T-EPIR schemes.*

**Definition 2.** *The secrecy rate is the amount of common randomness shared by the storage nodes relative to the file size, that is*

$$\rho_{\text{T-EPIR}} \triangleq \frac{H(S)}{H(W_k)}.$$

### III. Main Result

In this section, we summarize the main results of this paper.

**Theorem 1** (Capacity when $E \geq T$)**.** *For T-EPIR with $K$ files replicated at $N$ databases, where*

*any $T$ nodes may collude and an eavesdropper can tap in on the communication of any $E$ nodes, when $E \geq T$, the capacity is*

$$C_{\text{T-EPIR}} = \begin{cases} 1 - \frac{E}{N}, & \text{if } \rho_{\text{T-EPIR}} \geq \frac{E}{N-E} \\ 0, & \text{otherwise} \end{cases}.$$

*Remark:* For the detailed proof of Theorem 1, we refer to Section V and Section VI.B of our previous work [1].

**Theorem 2** (Outer Bound when $E \leq T$). *For T-EPIR with $K$ files replicated at $N$ databases, where any $T$ nodes may collude and an eavesdropper can tap in on the communication of any $E$ nodes, when $E \leq T$,*

$$R_{\text{T-EPIR}}^* \leq \overline{R}_{\text{T-EPIR}} = \left(1 - \frac{T}{N}\right) \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K}. \tag{3}$$

*The secrecy rate, i.e. the ratio of the amount of common randomness to the file size is at least* $\rho_{\text{T-EPIR}} \geq \frac{\frac{E}{N}\left(1-\left(\frac{T}{N}\right)^K\right)}{\left(1-\frac{T}{N}\right)\left(1-\frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}\right)}.$

*Remark:* The proof of the outer bound is in Section IV. The outer bound is tight, that is, it can be achieved and is hence the capacity of the problem for the four special cases below.

- Case 1 ($E = T$): From Theorem 1, the capacity is $C_{\text{T-EPIR}} = 1 - \frac{E}{N}$ when $E = T$. The outer bound in Theorem 2 is $\overline{R}_{\text{T-EPIR}} = \left(1 - \frac{T}{N}\right)\frac{1-\frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}}{1-\left(\frac{T}{N}\right)^K} = 1 - \frac{T}{N} = 1 - \frac{E}{N} = C_{\text{T-EPIR}}$ when $E = T$.

- Case 2 ($E = 0$): When there is no eavesdropper, *i.e.* $E = 0$, the problem reduce to the TPIR problem in [6], where the authors derive the capacity to be $C_{\text{TPIR}} = \frac{1-\frac{T}{N}}{1-\left(\frac{T}{N}\right)^K}$. The outer bound in Theorem 2 is $\overline{R}_{\text{T-EPIR}} = \left(1 - \frac{T}{N}\right)\frac{1-\frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}}{1-\left(\frac{T}{N}\right)^K} = \frac{1-\frac{T}{N}}{1-\left(\frac{T}{N}\right)^K} = C_{\text{TPIR}}$ when $E = 0$.

- Case 3 ($K \to \infty$): In our previous work [1], we derive the T-ESPIR capacity to be $C_{\text{T-ESPIR}} = 1 - \frac{\max(T,E)}{N} = 1 - \frac{T}{N}$ when $E \leq T$. As with all previous works for various scenarios of the PIR and SPIR problems, the PIR capacity reduces to the SPIR capacity when the number of files $K \to \infty$. The intuition is that, when the number of files increases, the penalty in the downloading rate to protect database-privacy for SPIR decays. When there are asymptotically infinitely many files, the information rate the user can learn about the

database from finite downloaded symbols vanishes. When the number of files $K$ tends to infinity, the outer bound tends to $\lim_{K \to \infty} \overline{R}_{\text{T-EPIR}} = \lim_{K \to \infty} \left(1 - \frac{T}{N}\right) \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K} = 1 - \frac{T}{N} = C_{\text{T-ESPIR}}$.

- Case 4 ($T = N$): When all databases collude, that is $T = N$, if furthermore $E = T = N$, the capacity is $0$ because the eavesdropper receives the same information as the user. If the user can decode $W_k$, so does the eavesdropper. Hence, the problem is non-trivial only if $E$ is strictly smaller than $T$. Suppose each file consists $L = N - E$ symbols from a large enough finite field $\mathbb{F}_q$, denoted by row vectors $W_k^{[1:L]}$ for $k \in [1 : K]$, consider the scheme below.

  The databases generate $KE$ uniformly i.i.d. symbols from $\mathbb{F}_q$, denoted by $K$ length-$E$ row vectors $S_k^{[1:E]}$ for $k \in [1 : K]$. Let $\mathbf{G}^{E \times N}$ be the generating matrix of an $(N, E)$-MDS code. The databases operate the $(N, E)$-MDS code on the common randomness vectors to obtain $K$ length-$N$ vectors $\bar{S}_k^{[1:N]} = S_k^{[1:E]} \mathbf{G}^{E \times N}$ for $k \in [1 : K]$, such that any $E$ symbols from $\bar{S}_k^{[1:N]}$ are uniformly identically distributed over $\mathbb{F}_q$. For each $k$, let $A_k^{[1:N]} = [\mathbf{0}^{[1 \times E]} W_k^{[1:L]}] + \bar{S}_k^{[1:N]}$ where $\mathbf{0}^{[1 \times E]}$ is a length-$E$ zero vector, the user downloads $A_k^n$ from database $n$ for each file index $k$. It can be checked that the user can decode $W_k$ (in fact the user can decode all files), and both user-privacy and system-privacy are guaranteed. The rate achieved by the scheme is $\frac{N-E}{NK}$.

  The outer bound in Theorem 2 is $\overline{R}_{\text{T-EPIR}} = \left(1 - \frac{T}{N}\right) \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K} = \frac{1 - \frac{E}{N} \cdot \left(\frac{T}{N}\right)^{K-1}}{1 + \frac{T}{N} + \cdots + \left(\frac{T}{N}\right)^{K-1}} = \frac{1 - \frac{E}{N}}{K} = \frac{N-E}{NK}$ when $T = N$, which is achieved by the scheme above.

**Theorem 3** (Inner Bound when $E \leq T$). *For T-EPIR with $K$ files replicated at $N$ databases, where any $T$ nodes may collude and an eavesdropper can tap in on the communication of any $E$ nodes, when $E \leq T$,*

$$R^*_{T\text{-}EPIR} \geq \underline{R}_{T\text{-}EPIR} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{E}{KN} \tag{4}$$

*Remark:* The inner bound is achieved by the scheme described in Section V. We discuss below the rate achieved by our scheme for the four special cases discussed above in which the outer bound in Theorem 2 is tight.

- Case 1 ($E = T$): The capacity of T-EPIR is $C_{\text{T-EPIR}} = 1 - \frac{E}{N} = 1 - \frac{T}{N}$ when $E = T$, that is,

(a) $N = 10, E = 3, K = 10$            (b) $N = 10, E = 3, K = 100$

Fig. 1: Plot of the bounds as functions of $\frac{T}{N}$.

the rate of $1 - \frac{T}{N}$ can be achieved by the scheme in our previous work [1]. When $E = T$, the rate achieved by the scheme in Section V is $\underline{R}_{\text{T-EPIR}} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{E}{KN} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{T}{KN} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} \cdot \left( 1 - \frac{1}{K} \left( \frac{T}{N} + (\frac{T}{N})^2 + \cdots + (\frac{T}{N})^K \right) \right)$, which is strictly smaller than $1 - \frac{T}{N}$ when $T \neq N$. Therefore, our scheme in Section V is not optimal when $E = T$. In other words, the inner bound in Theorem 3 is not tight for the case $E = T$.

- Case 2 ($E = 0$): When there is no eavesdropper hence $E = 0$, the rate achieved is $\underline{R}_{\text{T-EPIR}} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{E}{KN} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K}$, which matches with the TPIR capacity derived in [6], hence is optimal.

- Case 3 ($K \to \infty$): When the number of files $K$ tends to infinity, $\lim_{K \to \infty} \underline{R}_{\text{T-EPIR}} = \lim_{K \to \infty} \left( \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{E}{KN} \right) = 1 - \frac{T}{N}$. Hence, the inner bound tends to the T-ESPIR capacity as $K \to \infty$.

- Case 4 ($T = N$): When all databases collude hence $T = N$, the rate achieved by the scheme in this work is $\underline{R}_{\text{T-EPIR}} = \frac{1 - \frac{T}{N}}{1 - (\frac{T}{N})^K} - \frac{E}{KN} = \frac{1}{K} - \frac{E}{KN} = \frac{N - E}{KN}$, which matches the outer bound when $T = N$, hence is optimal.

In Figure 1 and Figure 2, the results of Theorems 1- 3 are plotted for several sets of parameters. It can be observed from the figures that when the number of messages $K$ increases, the gap between the inner and outer bounds decays and vanishes as $K \to \infty$.

(a) $N = 10, T = 7, K = 10$        (b) $N = 10, T = 7, K = 100$

Fig. 2: Plot of the bounds as functions of $\frac{E}{N}$.

## IV. OUTER BOUND WHEN $E \leq T$

In this section, we derive the outer bound presented in Theorem 2 for the PIR problem with $T$-colluding databases and $E$-eavesdropped databases when $E \leq T$. We start from the case when $K = 1$ and $K = 2$, then generalize to the case of arbitrary $K$ in Section IV-C.

### A. $K = 1$ Message

For any set of nodes $\mathcal{E} \subset [1 : N]$ with $|\mathcal{E}| = E$,

$$L = H(W_1) = H(W_1|\mathcal{Q}) - H(W_1|A^{[1]}_{[1:N]}, \mathcal{Q}) \tag{5}$$

$$= I(W_1; A^{[1]}_{[1:N]}|\mathcal{Q}) \tag{6}$$

$$= H(A^{[1]}_{[1:N]}|\mathcal{Q}) - H(A^{[1]}_{[1:N]}|W_1, \mathcal{Q}) \tag{7}$$

$$\leq H(A^{[1]}_{[1:N]}|\mathcal{Q}) - H(A^{[1]}_{\mathcal{E}}|W_1, \mathcal{Q}) \tag{8}$$

$$= H(A^{[1]}_{[1:N]}|\mathcal{Q}) - H(A^{[1]}_{\mathcal{E}}|\mathcal{Q}), \tag{9}$$

where (9) follows from system-privacy (2). Averaging over all $\mathcal{E}$ with size $E$ from $[1 : N]$, we have that

$$L \leq H(A^{[1]}_{[1:N]}|\mathcal{Q}) - \frac{1}{\binom{N}{E}} \sum_{\substack{\mathcal{E} \subset [1:N] \\ |\mathcal{E}| = E}} H(A^{[1]}_{\mathcal{E}}|\mathcal{Q}). \tag{10}$$

By Han's inequality [18],

$$\frac{1}{\binom{N}{E}} \sum_{\substack{\mathcal{E} \subset [1:N] \\ |\mathcal{E}| = E}} H(A_{\mathcal{E}}^{[1]}|\mathcal{Q}) \geq \frac{E}{N} H(A_{[1:N]}^{[1]}|\mathcal{Q}). \tag{11}$$

Therefore, $L \leq \left(1 - \frac{E}{N}\right) H(A_{[1:N]}^{[1]}|\mathcal{Q})$ and hence $R = \frac{L}{\sum_{n=1}^{N} H(A_n^{[1]})} \leq \frac{L}{H(A_{[1:N]}^{[1]}|\mathcal{Q})} \leq 1 - \frac{E}{N}$.

### B. $K = 2$ Messages

For any set of nodes $\mathcal{T} \subset [1:N]$ with $|\mathcal{T}| = T$, because of user-privacy, we can ignore the requested file index of $A_{\mathcal{T}}$,

$$L = H(W_1) = H(W_1) - H(W_1|A_{[1:N]}^{[1]}, \mathcal{Q}) \tag{12}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) - H(A_{[1:N]}^{[1]}|W_1, \mathcal{Q}) \tag{13}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) - H(A_{\mathcal{T}}|W_1, \mathcal{Q}) \tag{14}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) - \frac{T}{N} H(A_{[1:N]}^{[2]}|W_1, \mathcal{Q}), \tag{15}$$

where the last step (15) is obtained by averaging over all $\mathcal{T}$ with size $T$ and applying Han's inequality, similarly as (10) and (11) in Section IV-A. hence, we have that

$$H(A_{[1:N]}^{[2]}|W_1, \mathcal{Q}) \leq \frac{N}{T} \left(H(A_{[1:N]}^{[1]}|\mathcal{Q}) - L\right). \tag{16}$$

For any set of nodes $\mathcal{E} \subset [1:N]$ with $|\mathcal{E}| = E$, and any set of nodes $\mathcal{T} \subset [1:N]$ with $|\mathcal{T}| = T$,

$$2L = H(W_1, W_2) \tag{17}$$

$$= H(W_1, W_2|\mathcal{Q}) - H(W_1, W_2|A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}, \mathcal{Q}) \tag{18}$$

$$= I(W_1, W_2; A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}|\mathcal{Q}) \tag{19}$$

$$= H(A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}|\mathcal{Q}) - H(A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}|W_1, W_2, \mathcal{Q}) \tag{20}$$

$$\leq H(A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}|\mathcal{Q}) - H(A_{\mathcal{E}}|W_1, W_2, \mathcal{Q}) \tag{21}$$

$$= H(A_{[1:N]}^{[1]}, A_{[1:N]}^{[2]}|\mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{22}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{[1:N]}^{[2]}|A_{[1:N]}^{[1]}, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{23}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{[1:N]}^{[2]}|A_{[1:N]}^{[1]}, W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{24}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{[1:N]}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{25}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{[1:N]}^{[2]}|W_1, \mathcal{Q}) - H(A_{\mathcal{T}}|W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{26}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + \left(1 - \frac{T}{N}\right) H(A_{[1:N]}^{[2]}|W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{27}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + \left(1 - \frac{T}{N}\right) \frac{N}{T} \left(H(A_{[1:N]}^{[1]}|\mathcal{Q}) - L\right) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{28}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + \left(1 - \frac{T}{N}\right) \frac{N}{T} \left(H(A_{[1:N]}^{[1]}|\mathcal{Q}) - L\right) - \frac{E}{N} H(A_{[1:N]}^{[1]}|\mathcal{Q}) \tag{29}$$

$$= \left(\frac{N}{T} - \frac{E}{N}\right) H(A_{[1:N]}^{[1]}|\mathcal{Q}) - \left(\frac{N}{T} - 1\right) L, \tag{30}$$

where in (21) we can omit the message index because $\mathcal{E}$ is a set with size $E \leq T$. (22) follows from system-privacy (2). (24) is due to the fact that the user can decode $W_1$ from $A_{[1:N]}^{[1]}$ and $\mathcal{Q}$. (27) is obtained by averaging over all $\mathcal{T}$ with size $T$ and applying Han's inequality. (28) follows from (16). (29) is obtained by averaging over all $\mathcal{E}$ with size $E$ and applying Han's inequality.

Therefore, we have that $\left(\frac{N}{T} - \frac{E}{N}\right) H(A_{[1:N]}^{[1]}|\mathcal{Q}) \geq \left(\frac{N}{T} + 1\right) L$ and

$$R = \frac{L}{\sum_{n=1}^{N} H(A_n^{[1]})} \leq \frac{L}{H(A_{[1:N]}^{[1]}|\mathcal{Q})} \leq \frac{1 - \frac{E}{N} \cdot \frac{T}{N}}{1 + \frac{T}{N}}. \tag{31}$$

## C. $K \geq 3$ Messages

For any set of nodes $\mathcal{T} \subset [1:N]$ with $|\mathcal{T}| = T$, and its compliment set $\overline{\mathcal{T}} = [1:N] \setminus \mathcal{T}$, and for any $k \in [2:K]$,

$$H(A_{\overline{\mathcal{T}}}^{[k]}|A_{\mathcal{T}}, W_{[1:k-1]}, \mathcal{Q}) \tag{32}$$

$$= H(A_{[1:N]}^{[k]}|W_{[1:k-1]}, \mathcal{Q}) - H(A_{\mathcal{T}}|W_{[1:k-1]}, \mathcal{Q}) \tag{33}$$

$$\leq \left(1 - \frac{T}{N}\right) H(A_{[1:N]}^{[k]}|W_{[1:k-1]}, \mathcal{Q}), \tag{34}$$

where the last step follows by averaging over all $\mathcal{T}$ with size $T$ and applying Han's inequality.

From $A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]}$, the user can decode $W_{[1:k-1]}$, hence

$$(k-1)L = I(W_{[1:k-1]}; A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]}|\mathcal{Q}) \tag{35}$$

$$= H(A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]}|\mathcal{Q}) - H(A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]}|W_{[1:k-1]}, \mathcal{Q}) \tag{36}$$

$$\leq H(A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]} | \mathcal{Q}) - H(A_{\mathcal{T}} | W_{[1:k-1]}, \mathcal{Q}) \tag{37}$$

$$\leq H(A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]} | \mathcal{Q}) - \frac{T}{N} H(A_{[1:N]}^{[k]} | W_{[1:k-1]}, \mathcal{Q}), \tag{38}$$

where in (37), we can omit the message index of $A_{\mathcal{T}}$ because from user-privacy, the answers of any $T$ databases are independent of the message index. Similar as above, the last step follows by averaging over all $\mathcal{T}$ with size $T$ and applying Han's inequality. Because $A_{\mathcal{T}}$ is independent of the message index, we can set the index to $k$ in the last step.

Therefore, from (34) and (38), for any $k \in [2 : K]$,

$$H(A_{\mathcal{T}}^{[k]} | A_{\mathcal{T}}, W_{[1:k-1]}, \mathcal{Q}) \tag{39}$$

$$\leq \left(1 - \frac{T}{N}\right) H(A_{[1:N]}^{[k]} | W_{[1:k-1]}, \mathcal{Q}) \tag{40}$$

$$\leq \left(1 - \frac{T}{N}\right) \frac{N}{T} \left( H(A_{[1:N]}^{[1]}, \ldots, A_{[1:N]}^{[k-1]} | \mathcal{Q}) - (k-1)L \right) \tag{41}$$

$$= \left(\frac{N}{T} - 1\right) \left( H(A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, \ldots, A_{\mathcal{T}}^{[k-1]} | \mathcal{Q}) - (k-1)L \right) \tag{42}$$

$$= \left(\frac{N}{T} - 1\right) \left( H(A_{[1:N]}^{[1]} | \mathcal{Q}) + H(A_{\mathcal{T}}^{[2]}, \ldots, A_{\mathcal{T}}^{[k-1]} | A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, W_1, \mathcal{Q}) - (k-1)L \right) \tag{43}$$

$$\leq \left(\frac{N}{T} - 1\right) \left( H(A_{[1:N]}^{[1]} | \mathcal{Q}) + H(A_{\mathcal{T}}^{[2]}, \ldots, A_{\mathcal{T}}^{[k-1]} | A_{\mathcal{T}}, W_1, \mathcal{Q}) - (k-1)L \right) \tag{44}$$

$$\leq \left(\frac{N}{T} - 1\right) \left( H(A_{[1:N]}^{[1]} | \mathcal{Q}) + H(A_{\mathcal{T}}^{[2]} | A_{\mathcal{T}}, W_1, \mathcal{Q}) + \cdots + H(A_{\mathcal{T}}^{[k-1]} | A_{\mathcal{T}}, W_{[1:k-2]}, \mathcal{Q}) - (k-1)L \right), \tag{45}$$

where (43) holds because from $A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}$ and $\mathcal{Q}$ one can decode $W_1$. The last step is obtained by repeating the chain rule and by the fact that from $A_{\mathcal{T}}, A_{\mathcal{T}}^{[i]}$ and $\mathcal{Q}$ one can decode $W_i$ for $i = [2 : k - 2]$.

For any set of nodes $\mathcal{E} \subset [1 : N]$ with $|\mathcal{E}| = E$,

$$KL = H(W_{[1:K]}) \tag{46}$$

$$= I(W_{[1:K]}; A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, \ldots, A_{\mathcal{T}}^{[K]} | \mathcal{Q}) \tag{47}$$

$$= H(A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, \ldots, A_{\mathcal{T}}^{[K]} | \mathcal{Q}) - H(A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, \ldots, A_{\mathcal{T}}^{[K]} | W_{[1:K]}, \mathcal{Q}) \tag{48}$$

$$\leq H(A_{\mathcal{T}}, A_{\mathcal{T}}^{[1]}, \ldots, A_{\mathcal{T}}^{[K]} | \mathcal{Q}) - H(A_{\mathcal{E}} | W_{[1:K]}, \mathcal{Q}) \tag{49}$$

$$= H(A_{\mathcal{T}}, A_{\overline{\mathcal{T}}}^{[1]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|\mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{50}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{[1:N]}^{[1]}, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{51}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{[1:N]}^{[1]}, W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{52}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{53}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[3]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{\overline{\mathcal{T}}}^{[2]}, A_{\mathcal{T}}, W_1, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{54}$$

$$= H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[3]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{\overline{\mathcal{T}}}^{[2]}, A_{\mathcal{T}}, W_1, W_2, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{55}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[3]}, \ldots, A_{\overline{\mathcal{T}}}^{[K]}|A_{\mathcal{T}}, W_1, W_2, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{56}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[3]},|A_{\mathcal{T}}, W_1, W_2, \mathcal{Q}) + \ldots \tag{57}$$

$$+ H(A_{\overline{\mathcal{T}}}^{[K]},|A_{\mathcal{T}}, W_{[1:K-1]}, \mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{58}$$

$$\leq H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[3]},|A_{\mathcal{T}}, W_1, W_2, \mathcal{Q}) + \cdots + \left(\frac{N}{T} - 1\right)\left(H(A_{[1:N]}^{[1]}|\mathcal{Q}) + \right. \tag{59}$$

$$\left. H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + \cdots + H(A_{\overline{\mathcal{T}}}^{[K-1]}|A_{\mathcal{T}}, W_{[1:K-2]}, \mathcal{Q}) - (K-1)L\right) - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{60}$$

$$= \frac{N}{T}\left(H(A_{[1:N]}^{[1]}|\mathcal{Q}) + H(A_{\overline{\mathcal{T}}}^{[2]}|A_{\mathcal{T}}, W_1, \mathcal{Q}) + \cdots + H(A_{\overline{\mathcal{T}}}^{[K-1]}|A_{\mathcal{T}}, W_{[1:K-2]}, \mathcal{Q})\right) - \tag{61}$$

$$\left(\frac{N}{T} - 1\right)(K-1)L - H(A_{\mathcal{E}}|\mathcal{Q}) \tag{62}$$

$$\leq \left(\frac{N}{T}\right)^{K-1} H(A_{[1:N]}^{[1]}|\mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) - \left(1 - \frac{T}{N}\right)\left[\frac{N}{T}(K-1)L + \left(\frac{N}{T}\right)^2 (K-2)L\right. \tag{63}$$

$$\left. + \cdots + \left(\frac{N}{T}\right)^{K-1} L\right] \tag{64}$$

$$= \left(\frac{N}{T}\right)^{K-1} H(A_{[1:N]}^{[1]}|\mathcal{Q}) - H(A_{\mathcal{E}}|\mathcal{Q}) - \frac{\left(\frac{N}{T}\right)^K - \frac{N}{T}}{\frac{N}{T} - 1}L - (K-1)L \tag{65}$$

$$\leq \left(\left(\frac{N}{T}\right)^{K-1} - \frac{E}{N}\right) H(A_{[1:N]}^{[1]}|\mathcal{Q}) - \frac{\left(\frac{N}{T}\right)^K - \frac{N}{T}}{\frac{N}{T} - 1}L - (K-1)L, \tag{66}$$

where (50) is due to system-privacy (2). Steps (51)-(58) follows by repeating the chain rule and by the fact that from $A_{\mathcal{T}}$, $A_{\overline{\mathcal{T}}}^{[i]}$ and $\mathcal{Q}$ one can decode $W_i$ for $i = [1 : K-1]$. Step (60) follows by using inequality (45) for $k = K$. By iteratively using inequality (45) for $k = \{K - 1, K - 2, \ldots, 2\}$, we obtain (64). The last step follows by averaging over all $\mathcal{E}$ with size $E$ and applying Han's

inequality.

Therefore, $\left( \left(\frac{N}{T}\right)^{K-1} - \frac{E}{N}\right) H(A^{[1]}_{[1:N]}|\mathcal{Q}) \geq \frac{\left(\frac{N}{T}\right)^K - 1}{\frac{N}{T} - 1} L$, and hence

$$R = \frac{L}{\sum_{n=1}^N H(A_n^{[1]})} \leq \frac{L}{H(A^{[1]}_{[1:N]}|\mathcal{Q})} \leq \left(\frac{N}{T} - 1\right)\frac{\left(\frac{N}{T}\right)^{K-1} - \frac{E}{N}}{\left(\frac{N}{T}\right)^K - 1} \tag{67}$$

$$= \left(1 - \frac{T}{N}\right)\frac{1 - \frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}}{1 - \left(\frac{T}{N}\right)^K} \tag{68}$$

$$= \overline{R}_{\text{T-EPIR}}. \tag{69}$$

To obtain a lower bound on the amount of common randomness needed to guarantee system-privacy, for any set of nodes $\mathcal{E} \subset [1:N]$ with size $|\mathcal{E}| = E$,

$$0 = I(A_\mathcal{E}; W_{[1:K]}|\mathcal{Q}) \tag{70}$$

$$= H(A_\mathcal{E}|\mathcal{Q}) - H(A_\mathcal{E}|W_{[1:K]}, \mathcal{Q}) \tag{71}$$

$$= H(A_\mathcal{E}|\mathcal{Q}) - H(A_\mathcal{E}|W_{[1:K]}, \mathcal{Q}) + H(A_\mathcal{E}|W_{[1:K]}, S, \mathcal{Q}) \tag{72}$$

$$= H(A_\mathcal{E}|\mathcal{Q}) - I(S; A_\mathcal{E}|W_{[1:K]}, \mathcal{Q}) \tag{73}$$

$$= H(A_\mathcal{E}|\mathcal{Q}) - H(S|W_{[1:K]}, \mathcal{Q}) + H(S|A_\mathcal{E}, W_{[1:K]}, \mathcal{Q}) \tag{74}$$

$$\geq H(A_\mathcal{E}|\mathcal{Q}) - H(S), \tag{75}$$

where (72) holds because $A_\mathcal{E}$ is a deterministic function of $W_{[1:K]}, S$ and $\mathcal{Q}$. By averaging over all $\mathcal{E}$ with size $E$ and applying Han's inequality,

$$H(S) \geq \frac{1}{\binom{N}{E}}\sum_{\substack{\mathcal{E}\subset[1:N]\\|\mathcal{E}|=E}} H(A_\mathcal{E}|\mathcal{Q}) \geq \frac{E}{N}H(A^{[1]}_{[1:N]}|\mathcal{Q}) \tag{76}$$

$$\geq \frac{\frac{E}{N}\left(1 - \left(\frac{T}{N}\right)^K\right)}{\left(1 - \frac{T}{N}\right)\left(1 - \frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}\right)} L. \tag{77}$$

Therefore, $\rho_{\text{T-EPIR}} = \frac{H(S)}{L} \geq \frac{\frac{E}{N}\left(1-\left(\frac{T}{N}\right)^K\right)}{\left(1-\frac{T}{N}\right)\left(1-\frac{E}{N}\cdot\left(\frac{T}{N}\right)^{K-1}\right)}.$

## V. Inner bound when $E \leq T$

In this section, we present an achievable scheme for the case when the eavesdropper can tap in on any $E$ databases where $E \leq T$. The scheme is modified from the TPIR scheme in [6], by downloading $K$ rounds where each round use the scheme in [6] with different part of the files and different part of the common randomness generated by the databases. The three principles in [6] still apply in our scheme.

1) Symmetry across databases

2) Symmetry of file indices within the queries to each database

3) Exploiting the side information of undesired files to retrieve the desired file information

Specifically, the new ingredient of our scheme lies in iterating the scheme in $K$ rounds to ensure each file is mixed with the common randomness in the same way, hence to fulfill principle 2. In the following, we firstly introduce five examples. We explain in details of the examples in Section V-A and Section V-E about the decodability of the scheme, the guarantee of user-privacy and system-privacy, and only show the construction of the other three examples. Finally in Section V-F, we show the scheme for general parameters of $N, K, T, E$.

We first reprise the following lemma from [6]. The lemma states that by multiplying deterministic full rank matrices on uniformly i.i.d. random matrices, the statistics of the random matrices remain unchanged. The proof can be found in [6].

**Lemma 4** ( [6]). *Let $\mathbf{S}_1, \mathbf{S}_2, \ldots, \mathbf{S}_K \in \mathbb{F}_q^{\alpha \times \alpha}$ be $K$ random matrices, drawn independently and uniformly from all $\alpha \times \alpha$ full-rank matrices over $\mathbb{F}_q$. Let $\mathbf{G}_1, \mathbf{G}_2, \ldots, \mathbf{G}_K \mathbb{F}_q^{\beta \times \beta}$ be $K$ invertible square matrices of dimension $\beta \times \beta$ over $\mathbb{F}_q$ where $\beta \leq \alpha$. Let $\mathcal{I}_1, \mathcal{I}_2, \ldots, \mathcal{I}_K \in \mathbb{N}^{1 \times \beta}$ be $K$ index vectors, each containing $\beta$ distinct indices from $[1 : \alpha]$, then*

$$(\mathbf{S}_1[:, \mathcal{I}_1]\mathbf{G}_1, \mathbf{S}_2[:, \mathcal{I}_2]\mathbf{G}_2, \ldots, \mathbf{S}_K[:, \mathcal{I}_K]\mathbf{G}_K) \sim (\mathbf{S}_1[:, (1 : \beta)], \mathbf{S}_2[:, (1 : \beta)], \ldots, \mathbf{S}_K[:, (1 : \beta)]) \tag{78}$$

*where $\mathbf{S}_i[:, \mathcal{I}_i]$ denotes the $\alpha \times \beta$ matrix comprised of the columns of $\mathbf{S}_i$ with indices in $\mathcal{I}_i$, and $\sim$ denotes the relation that the random variables on both sides are identically distributed.*

| DB1 | DB2 | DB3 |
|---|---|---|
| $a_1^{(r)}, a_2^{(r)}$ | $a_3^{(r)}, a_4^{(r)}$ | $a_5^{(r)}, a_6^{(r)}$ |
| $b_1^{(r)}, b_2^{(r)}$ | $b_3^{(r)}, b_4^{(r)}$ | $b_5^{(r)}, b_6^{(r)}$ |
| $a_7^{(r)} + b_7^{(r)}$ | $a_8^{(r)} + b_8^{(r)}$ | $a_9^{(r)} + b_9^{(r)}$ |

TABLE I: The download scheme for each round $r$, where $r = 1$ and $r = 2$.

*A. Example: $N = 3$ databases, $K = 2$ files, $T = 2$ colluding databases, $E = 1$ eavesdropped database*

Suppose each file contains $L = 13$ symbols from a sufficiently large finite field $\mathbb{F}_q$, $W_1 = W_1^{[1:13]}$ and $W_2 = W_2^{[1:13]}$ are represented as length-13 vectors over $\mathbb{F}_q$. W.l.o.g., assume the user wants to retrieve $W_1$.

The user downloads in two rounds, with 15 symbols in each round as described in Table I and with detailed formulation below. The databases generate 10 uniformly random symbols, 5 for each round, denoted as $(S_{[1:5]}^{(1)}, S_{[1:5]}^{(2)})$. The scheme achieves the rate $R = 13/30$.

Let $\{\lambda_1, \ldots, \lambda_9\}$ be 9 distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}_{[1:7]}^{7 \times 9}$ and $\mathbf{G}_{[8:9]}^{2 \times 9}$ be two generating matrices of MDS codes as follows,

$$\mathbf{G}_{[1:7]}^{7 \times 9} = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_9 \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^6 & \lambda_2^6 & \ldots & \lambda_9^6 \end{bmatrix}, \tag{79}$$

$$\mathbf{G}_{[8:9]}^{2 \times 9} = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_9 \end{bmatrix} \cdot diag(\lambda_1^7, \lambda_2^7, \ldots, \lambda_9^7) \tag{80}$$

$$= \begin{bmatrix} \lambda_1^7 & \lambda_2^7 & \ldots & \lambda_9^7 \\ \lambda_1^8 & \lambda_2^8 & \ldots & \lambda_9^8 \end{bmatrix}. \tag{81}$$

Let $\mathbf{G} = [\mathbf{G}_{[1:7]}^{7 \times 9} \ \mathbf{G}_{[8:9]}^{2 \times 9}]^{\mathrm{T}}$, then $\mathbf{G}$ is a $9 \times 9$ invertible matrix. Similarly, let $\mathbf{G}_{[1:6]}^{6 \times 9}$ and $\mathbf{G}_{[7:9]}^{3 \times 9}$ be composed of the first six rows and the last three rows of $\mathbf{G}$ respectively.

The user privately generates matrices $\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3, \mathbf{S}_4 \in \mathbb{F}_q^{9 \times 9}$ uniformly and independently from all $9 \times 9$ invertible matrices over $\mathbb{F}_q$.

Let $\mathbf{G}_1^{6 \times 9}$ be the generating matrix of a $(9, 6)$-MDS code.

*Round 1:*

$$a_{[1:9]}^{(1)} = \left( W_1^{[1:7]} \mathbf{G}_{[1:7]}^{7\times9} + [S_1^{(1)} S_2^{(1)}] \mathbf{G}_{[8:9]}^{2\times9} \right) \mathbf{S}_1 \tag{82}$$

$$b_{[1:9]}^{(1)} = \left( W_2^{[8:13]} \mathbf{G}_{[1:6]}^{6\times9} + [S_3^{(1)} S_4^{(1)} S_5^{(1)}] \mathbf{G}_{[7:9]}^{3\times9} \right) \mathbf{S}_2[:, (1:6)] \mathbf{G}_1^{6\times9} \tag{83}$$

*Round 2:*

$$a_{[1:9]}^{(2)} = \left( W_1^{[8:13]} \mathbf{G}_{[1:6]}^{6\times9} + [S_3^{(2)} S_4^{(2)} S_5^{(2)}] \mathbf{G}_{[7:9]}^{3\times9} \right) \mathbf{S}_3 \tag{84}$$

$$b_{[1:9]}^{(2)} = \left( W_2^{[1:7]} \mathbf{G}_{[1:7]}^{7\times9} + [S_1^{(2)} S_2^{(2)}] \mathbf{G}_{[8:9]}^{2\times9} \right) \mathbf{S}_4[:, (1:6)] \mathbf{G}_1^{6\times9} \tag{85}$$

**Correctness:** In round 1, the user can solve $b_{[7:9]}^{(1)}$ from $b_{[1:6]}^{(1)}$, because $\mathbf{G}_3^{6\times9}$ is the generating matrix of a $(9,6)$-MDS code. Therefore, the user can cancel the interference $b_{[7:9]}^{(1)}$ and obtain $a_{[7:9]}^{(1)}$. From $a_{[1:9]}^{(1)}$, the user can solve $W_1^{[1:7]}$, because $a_{[1:9]}^{(1)} = [W_1^{[1:7]} S_1^{(1)} S_2^{(1)}] \mathbf{GS}_1$, where $\mathbf{G}$ and $\mathbf{S}_1$ are invertible matrices. Similarly in round 2, the user can solve $W_1^{[8:13]}$. Hence, the user can solve all 13 symbols of $W_1$.

**User-privacy:** Any $T = 2$ databases may collude and observe the queries composed of 6 symbols from $a_{[1:9]}^{(r)}$ and $b_{[1:9]}^{(r)}$ for each round. Let $\mathcal{I}_a, \mathcal{I}_b$ denote the indices of the symbols observed by the colluding databases,

$$\left( a_{\mathcal{I}_a}^{(1)}, a_{\mathcal{I}_a}^{(2)}, b_{\mathcal{I}_b}^{(1)}, b_{\mathcal{I}_b}^{(2)} \right) \tag{86}$$

$$= \left( \left[ W_1^{[1:7]} S_1^{(1)} S_2^{(1)} \right] \mathbf{GS}_1[:, \mathcal{I}_a], \left[ W_1^{[8:13]} S_3^{(2)} S_4^{(2)} S_5^{(2)} \right] \mathbf{GS}_3[:, \mathcal{I}_a], \tag{87}$$

$$\left[ W_2^{[8:13]} S_3^{(1)} S_4^{(1)} S_5^{(1)} \right] \mathbf{GS}_2[:, (1:6)] \mathbf{G}_1^{6\times9}[:, \mathcal{I}_b], \left[ W_2^{[1:7]} S_1^{(2)} S_2^{(2)} \right] \mathbf{GS}_4[:, (1:6)] \mathbf{G}_1^{6\times9}[:, \mathcal{I}_b] \right) \tag{88}$$

$$\sim \left( \left[ W_1^{[1:7]} S_1^{(1)} S_2^{(1)} \right] \mathbf{S}_1[:, (1:6)], \left[ W_1^{[8:13]} S_3^{(2)} S_4^{(2)} S_5^{(2)} \right] \mathbf{S}_3[:, (1:6)], \tag{89}$$

$$\left[ W_2^{[8:13]} S_3^{(1)} S_4^{(1)} S_5^{(1)} \right] \mathbf{S}_2[:, (1:6)], \left[ W_2^{[1:7]} S_1^{(2)} S_2^{(2)} \right] \mathbf{S}_4[:, (1:6)] \right). \tag{90}$$

The two rounds of download can be randomized by the user. Therefore, the symbols observed by the two databases are obtained by random mappings from linear combinations of $W_1$ and $W_2$ and the random symbols $S_{[1:5]}^{(1)}, S_{[1:5]}^{(2)}$ generated by the databases in the same way, where the randomness of the mapping is privately generated by the user and unavailable to the databases. Hence, user-privacy is guaranteed.

**System-privacy:** The eavesdropper can tap in on an arbitrary database. Because the scheme

is symmetric across the databases, w.l.o.g., assume DB1 is eavesdropped. In round 1, from equation (82) $a_1^{(1)}, a_2^{(1)}$ are constructed by adding linearly independent combinations of $S_1^{(1)}, S_2^{(1)}$. Similarly from equation (83), $b_1^{(1)}, b_2^{(1)}, b_7^{(1)}$ are constructed by adding linearly independent combinations of $S_3^{(1)}, S_4^{(1)}, S_5^{(1)}$. Specifically, denote the five answers from DB1 in round 1 by $A_{DB1}^{(1)}$, the linear combinations of the $S_{[1:5]}^{(1)}$ added to the answers are constructed by,

$$[S_1^{(1)} S_2^{(1)} S_3^{(1)} S_4^{(1)} S_5^{(1)}] \cdot \begin{bmatrix} \left[\mathbf{G}_{[8:9]}^{2\times9}\mathbf{S}_1[:,(1:2)]\right]^{2\times2} & \mathbf{0}^{2\times2} & \left[\mathbf{G}_{[8:9]}^{2\times9}\mathbf{S}_1[:,7]\right]^{2\times1} \\ \mathbf{0}^{3\times2} & \left[\mathbf{G}_{[7:9]}^{2\times9}\mathbf{S}_2[:,(1:6)]\mathbf{G}_1^{6\times9}[:,(1,2,7)]\right]^{3\times3} \end{bmatrix} \cdot \tag{91}$$

It can be checked that the $5 \times 5$ matrix in (91) is invertible. Therefore, $H(A_{DB1}^{(1)}) = H(A_{DB1}^{(1)}|W_1, W_2) = 5\log q$. Hence, $I(A_{DB1}^{(1)}; W_1, W_2) = 0$. The construction of symbols for round 2 are in a similar way, by adding linearly independent combinations of $S_{[1:5]}^{(2)}$. Because the 10 symbols $S_{[1:5]}^{(1)}, S_{[1:5]}^{(2)}$ are independently and uniformly chosen from $\mathbb{F}_q$, we have $I(A_{DB1}^{(1)}, A_{DB1}^{(2)}; W_1, W_2) = 0$ and hence the eavesdropper obtains no information regarding the database $W_1, W_2$.

*B. Example: $N = 4$ databases, $K = 2$ files, $T = 2$ colluding databases, $E = 1$ eavesdropped database*

Suppose each file consists of $L = 13$ symbols and is represented as a length-13 row vector over a sufficiently large field $\mathbb{F}_q$, denoted by $W_1 = W_1^{[1:13]}$ and $W_2 = W_2^{[1:13]}$. The user downloads two rounds. For each round, the user downloads 12 symbols. The databases generate 6 uniformly random symbols $S_{[1:3]}^{(1)}, S_{[1:3]}^{(2)}$. The scheme achieves the rate $R = 13/24$.

The user privately generates matrices $\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3, \mathbf{S}_4 \in \mathbb{F}_q^{8\times8}$ uniformly and independently from all $8 \times 8$ invertible matrices over $\mathbb{F}_q$.

Let $\{\lambda_1, \ldots, \lambda_8\}$ be 8 distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}$ be a $8 \times 8$ matrix defined as follows,

$$\mathbf{G} = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_8 \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^7 & \lambda_2^7 & \ldots & \lambda_8^7 \end{bmatrix}, \tag{92}$$

it is direct that $\mathbf{G}$ is an invertible matrix. Let $\mathbf{G}_{[1:7]}^{7\times8}$ and $\mathbf{G}_{[8]}^{1\times8}$ be matrices composed of the first 7 rows and the 8th row respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:7]}^{7\times8} \ \mathbf{G}_{[8]}^{1\times8}]^{\mathrm{T}}$. Similarly, let $\mathbf{G}_{[1:6]}^{6\times8}$

and $\mathbf{G}_{[7:8]}^{2\times 8}$ be matrices composed of the first 6 rows and the last 2 rows respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:6]}^{6\times 8}\ \mathbf{G}_{[7:8]}^{2\times 8}]^{\mathrm{T}}$. The matrices $\mathbf{G}_{[1:7]}^{7\times 8}$, $\mathbf{G}_{[8]}^{1\times 8}$, $\mathbf{G}_{[1:6]}^{6\times 8}$ and $\mathbf{G}_{[7:8]}^{2\times 8}$ are generating matrices of MDS codes with corresponding dimensions.

Let $\mathbf{G}_1^{4\times 8}$ be the generating matrix of a $(8,4)$-MDS code.

| DB1 | DB2 | DB3 | DB4 |
|:---:|:---:|:---:|:---:|
| $a_1^{(r)}$ | $a_2^{(r)}$ | $a_3^{(r)}$ | $a_4^{(r)}$ |
| $b_1^{(r)}$ | $b_2^{(r)}$ | $b_3^{(r)}$ | $b_4^{(r)}$ |
| $a_5^{(r)} + b_5^{(r)}$ | $a_6^{(r)} + b_6^{(r)}$ | $a_7^{(r)} + b_7^{(r)}$ | $a_8^{(r)} + b_8^{(r)}$ |

*Round 1:*

$$a_{[1:8]}^{(1)} = \left( W_1^{[1:7]}\mathbf{G}_{[1:7]}^{7\times 8} + S_1^{(1)}\mathbf{G}_{[8]}^{1\times 8} \right) \mathbf{S}_1 \tag{93}$$

$$b_{[1:8]}^{(1)} = \left( W_2^{[8:13]}\mathbf{G}_{[1:6]}^{6\times 8} + [S_2^{(1)} S_3^{(1)}]\mathbf{G}_{[7:8]}^{2\times 8} \right) \mathbf{S}_2[:,(1:4)]\mathbf{G}_1^{4\times 8} \tag{94}$$

*Round 2:*

$$a_{[1:8]}^{(2)} = \left( W_1^{[8:13]}\mathbf{G}_{[1:6]}^{6\times 8} + [S_2^{(2)} S_3^{(2)}]\mathbf{G}_{[7:8]}^{2\times 8} \right) \mathbf{S}_3 \tag{95}$$

$$b_{[1:8]}^{(2)} = \left( W_2^{[1:7]}\mathbf{G}_{[1:7]}^{7\times 8} + S_1^{(2)}\mathbf{G}_{[8]}^{1\times 8} \right) \mathbf{S}_4[:,(1:4)]\mathbf{G}_1^{4\times 8} \tag{96}$$

*C. Example: $N = 4$ databases, $K = 2$ files, $T = 3$ colluding databases, $E = 1$ eavesdropped database*

Suppose each file consists of $L = 25$ symbols and is represented as a length-25 row vector over a sufficiently large field $\mathbb{F}_q$, denoted by $W_1 = W_1^{[1:25]}$ and $W_2 = W_2^{[1:25]}$. The user downloads two rounds. For each round, the user downloads 28 symbols. The databases generate 14 uniformly random symbols $S_{[1:7]}^{(1)}, S_{[1:7]}^{(2)}$. The scheme achieves the rate $R = 25/56$.

The user privately generates matrices $\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3, \mathbf{S}_4 \in \mathbb{F}_q^{16\times 16}$ uniformly and independently from all $16 \times 16$ invertible matrices over $\mathbb{F}_q$.

Let $\{\lambda_1, \ldots, \lambda_{16}\}$ be 16 distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}$ be a $16 \times 16$ matrix defined as follows,

$$\mathbf{G}_2 = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_{16} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{15} & \lambda_2^{15} & \ldots & \lambda_{16}^{15} \end{bmatrix}, \tag{97}$$

it is direct that $\mathbf{G}$ is an invertible matrix. Let $\mathbf{G}_{[1:13]}^{13\times16}$ and $\mathbf{G}_{[14:16]}^{3\times16}$ be matrices composed of the first 13 rows and the last 3 rows respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:13]}^{13\times16}\ \mathbf{G}_{[14:16]}^{3\times16}]^{\mathsf{T}}$. Similarly, let $\mathbf{G}_{[1:12]}^{12\times16}$ and $\mathbf{G}_{[13:16]}^{4\times16}$ be matrices composed of the first 12 rows and the last 4 rows respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:12]}^{12\times16}\ \mathbf{G}_{[13:16]}^{4\times16}]^{\mathsf{T}}$. The matrices $\mathbf{G}_{[1:13]}^{13\times16}$, $\mathbf{G}_{[14:16]}^{3\times16}$, $\mathbf{G}_{[1:12]}^{12\times16}$ and $\mathbf{G}_{[13:16]}^{4\times16}$ are generating matrices of MDS codes with corresponding dimensions.

Let $\mathbf{G}_1^{12\times16}$ be the generating matrix of a $(16, 12)$-MDS code.

| DB1 | DB2 | DB3 | DB4 |
|---|---|---|---|
| $a_1^{(r)}, a_2^{(r)}, a_3^{(r)}$ | $a_4^{(r)}, a_5^{(r)}, a_6^{(r)}$ | $a_7^{(r)}, a_8^{(r)}, a_9^{(r)}$ | $a_{10}^{(r)}, a_{11}^{(r)}, a_{12}^{(r)}$ |
| $b_1^{(r)}, b_2^{(r)}, b_3^{(r)}$ | $b_4^{(r)}, b_5^{(r)}, b_6^{(r)}$ | $b_7^{(r)}, b_8^{(r)}, b_9^{(r)}$ | $b_{10}^{(r)}, b_{11}^{(r)}, b_{12}^{(r)}$ |
| $a_{13}^{(r)} + b_{13}^{(r)}$ | $a_{14}^{(r)} + b_{14}^{(r)}$ | $a_{15}^{(r)} + b_{15}^{(r)}$ | $a_{16}^{(r)} + b_{16}^{(r)}$ |

*Round 1:*

$$a_{[1:16]}^{(1)} = \left( W_1^{[1:13]}\mathbf{G}_{[1:13]}^{13\times16} + [S_1^{(1)}S_2^{(1)}S_3^{(1)}]\mathbf{G}_{[14:16]}^{3\times16} \right)\mathbf{S}_1 \tag{98}$$

$$b_{[1:16]}^{(1)} = \left( W_2^{[14:25]}\mathbf{G}_{[1:12]}^{12\times16} + [S_4^{(1)}S_5^{(1)}S_6^{(1)}S_7^{(1)}]\mathbf{G}_{[13:16]}^{4\times16} \right)\mathbf{S}_2[:, (1:12)]\mathbf{G}_1^{12\times16} \tag{99}$$

*Round 2:*

$$a_{[1:16]}^{(2)} = \left( W_1^{14:25]}\mathbf{G}_{[1:12]}^{12\times16} + [S_4^{(2)}S_5^{(2)}S_6^{(2)}S_7^{(2)}]\mathbf{G}_{[13:16]}^{4\times16} \right)\mathbf{S}_3 \tag{100}$$

$$b_{[1:16]}^{(2)} = \left( W_2^{[1:13]}\mathbf{G}_{[1:13]}^{13\times16} + [S_1^{(2)}S_2^{(2)}S_3^{(2)}]\mathbf{G}_{[14:16]}^{3\times16} \right)\mathbf{S}_4[:, (1:12)]\mathbf{G}_1^{12\times16} \tag{101}$$

*D. Example: $N = 4$ databases, $K = 2$ files, $T = 3$ colluding databases, $E = 2$ eavesdropped databases*

Suppose each file contains $L = 18$ symbols and is represented as a length-18 row vector over a sufficiently large field $\mathbb{F}_q$, denoted by $W_1 = W_1^{[1:18]}$ and $W_2 = W_2^{[1:18]}$. The user downloads two rounds. For each round, the user downloads 28 symbols. The databases generate 28 uniformly random symbols $S_{[1:14]}^{(1)}, S_{[1:14]}^{(2)}$. The scheme achieves the rate $R = 18/56 = 9/28$.

The user privately generates matrices $\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3, \mathbf{S}_4 \in \mathbb{F}_q^{16\times16}$ uniformly and independently from all $16 \times 16$ invertible matrices over $\mathbb{F}_q$.

Let $\{\lambda_1, \ldots, \lambda_{16}\}$ be 16 distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}$ be a $16 \times 16$ matrix defined

as follows,

$$\mathbf{G} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \lambda_1 & \lambda_2 & \dots & \lambda_{16} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{15} & \lambda_2^{15} & \dots & \lambda_{16}^{15} \end{bmatrix}, \tag{102}$$

it is direct that $\mathbf{G}$ is an invertible matrix. Let $\mathbf{G}_{[1:10]}^{10\times16}$ and $\mathbf{G}_{[11:16]}^{6\times16}$ be matrices composed of the first 10 rows and the last 6 rows respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:10]}^{10\times16} \ \mathbf{G}_{[11:16]}^{6\times16}]^{\mathrm{T}}$. Similarly, let $\mathbf{G}_{[1:8]}^{8\times16}$ and $\mathbf{G}_{[9:16]}^{8\times16}$ be matrices composed of the first 8 rows and the last 8 rows respectively, such that $\mathbf{G} = [\mathbf{G}_{[1:8]}^{8\times16} \ \mathbf{G}_{[9:16]}^{8\times16}]^{\mathrm{T}}$. The matrices $\mathbf{G}_{[1:10]}^{10\times16}$, $\mathbf{G}_{[11:16]}^{6\times16}$, $\mathbf{G}_{[1:8]}^{8\times16}$ and $\mathbf{G}_{[9:16]}^{8\times16}$ are generating matrices of MDS codes with corresponding dimensions.

Let $\mathbf{G}_1^{12\times16}$ be the generating matrix of a $(16, 12)$-MDS code.

| DB1 | DB2 | DB3 | DB4 |
|---|---|---|---|
| $a_1^{(r)}, a_2^{(r)}, a_3^{(r)}$ | $a_4^{(r)}, a_5^{(r)}, a_6^{(r)}$ | $a_7^{(r)}, a_8^{(r)}, a_9^{(r)}$ | $a_{10}^{(r)}, a_{11}^{(r)}, a_{12}^{(r)}$ |
| $b_1^{(r)}, b_2^{(r)}, b_3^{(r)}$ | $b_4^{(r)}, b_5^{(r)}, b_6^{(r)}$ | $b_7^{(r)}, b_8^{(r)}, b_9^{(r)}$ | $b_{10}^{(r)}, b_{11}^{(r)}, b_{12}^{(r)}$ |
| $a_{13}^{(r)} + b_{13}^{(r)}$ | $a_{14}^{(r)} + b_{14}^{(r)}$ | $a_{15}^{(r)} + b_{15}^{(r)}$ | $a_{16}^{(r)} + b_{16}^{(r)}$ |

*Round 1:*

$$a_{[1:16]}^{(1)} = \left( W_1^{[1:10]} \mathbf{G}_{[1:10]}^{10\times16} + S_{[1:6]}^{(1)} \mathbf{G}_{[11:16]}^{6\times16} \right) \mathbf{S}_1 \tag{103}$$

$$b_{[1:16]}^{(1)} = \left( W_2^{[11:18]} \mathbf{G}_{[1:8]}^{8\times16} + S_{[7:14]}^{(1)} \mathbf{G}_{[9:16]}^{8\times16} \right) \mathbf{S}_2[:, (1:12)] \mathbf{G}_1^{12\times16} \tag{104}$$

*Round 2:*

$$a_{[1:16]}^{(2)} = \left( W_1^{[11:18]} \mathbf{G}_{[1:8]}^{8\times16} + S_{[7:14]}^{(2)} \mathbf{G}_{[9:16]}^{8\times16} \right) \mathbf{S}_3 \tag{105}$$

$$b_{[1:16]}^{(2)} = \left( W_2^{[1:10]} \mathbf{G}_{[1:10]}^{10\times16} + S_{[1:6]}^{(2)} \mathbf{G}_{[11:16]}^{6\times16} \right) \mathbf{S}_4[:, (1:12)] \mathbf{G}_1^{12\times16} \tag{106}$$

*E. Example: $N = 3$ databases, $K = 3$ files, $T = 2$ colluding databases, $E = 1$ eavesdropped database*

Suppose each file contains $L = 62$ symbols. Let the symbols of each file be randomly permuted (the randomness is generated privately by the user) and be represented as a length-62 row vector over a sufficiently large field $\mathbb{F}_q$, denoted by $W_1 = W_1^{[1:62]}$, $W_2 = W_2^{[1:62]}$ and $W_3 = W_3^{[1:62]}$. The user downloads three rounds. For each round, the user downloads 57 symbols. The databases

generate 57 uniformly random symbols, 19 for each round and denoted as $S^{(1)}_{[1:19]}, S^{(2)}_{[1:19]}, S^{(3)}_{[1:19]}$, for protecting the database from the eavesdropper. The scheme achieves the rate $R = 62/171$.

Let $\{\lambda_1, \ldots, \lambda_{27}\}$ be 27 distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}$ be a $27 \times 27$ matrix defined as follows,

$$\mathbf{G} = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_{27} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{26} & \lambda_2^{26} & \ldots & \lambda_{27}^{26} \end{bmatrix}, \tag{107}$$

it is direct that $\mathbf{G}$ is an invertible matrix. Let $\mathbf{G}^{18 \times 27}_{[1:18]}$ and $\mathbf{G}^{9 \times 27}_{[19:27]}$ be matrices composed of the first 18 rows and the last 9 rows respectively, such that $\mathbf{G} = [\mathbf{G}^{18 \times 27}_{[1:18]} \ \mathbf{G}^{9 \times 27}_{[19:27]}]^{\mathrm{T}}$. Similarly, let $\mathbf{G}^{21 \times 27}_{[1:21]}$ and $\mathbf{G}^{6 \times 27}_{[22:27]}$ be matrices composed of the first 21 rows and the last 6 rows respectively, and $\mathbf{G}^{23 \times 27}_{[1:23]}$ and $\mathbf{G}^{4 \times 27}_{[24:27]}$ be matrices composed of the first 23 rows and the last 4 rows respectively. The matrices $\mathbf{G}^{18 \times 27}_{[1:18]}$, $\mathbf{G}^{9 \times 27}_{[19:27]}$, $\mathbf{G}^{21 \times 27}_{[1:21]}$, $\mathbf{G}^{6 \times 27}_{[22:27]}$, $\mathbf{G}^{23 \times 27}_{[1:23]}$ and $\mathbf{G}^{4 \times 27}_{[24:27]}$ are generating matrices of MDS codes with corresponding dimensions.

The user privately generates 9 matrices $\mathbf{S}^{(1)}_{[1:3]}, \mathbf{S}^{(2)}_{[1:3]}, \mathbf{S}^{(3)}_{[1:3]} \in \mathbb{F}_q^{27 \times 27}$ uniformly and independently from all $27 \times 27$ invertible matrices over $\mathbb{F}_q$.

Let $\mathbf{G}_1^{12 \times 18}$ and $\mathbf{G}_2^{6 \times 9}$ be the generating matrices of a $(18, 12)$-MDS code and a $(9, 6)$-MDS code respectively.

| DB1 | DB2 | DB3 |
|---|---|---|
| $a_1^{(r)}, a_2^{(r)}, a_3^{(r)}, a_4^{(r)}$ | $a_5^{(r)}, a_6^{(r)}, a_7^{(r)}, a_8^{(r)}$ | $a_9^{(r)}, a_{10}^{(r)}, a_{11}^{(r)}, a_{12}^{(r)}$ |
| $b_1^{(r)}, b_2^{(r)}, b_3^{(r)}, b_4^{(r)}$ | $b_5^{(r)}, b_6^{(r)}, b_7^{(r)}, b_8^{(r)}$ | $b_9^{(r)}, b_{10}^{(r)}, b_{11}^{(r)}, b_{12}^{(r)}$ |
| $c_1^{(r)}, c_2^{(r)}, c_3^{(r)}, c_4^{(r)}$ | $c_5^{(r)}, c_6^{(r)}, c_7^{(r)}, c_8^{(r)}$ | $c_9^{(r)}, c_{10}^{(r)}, c_{11}^{(r)}, c_{12}^{(r)}$ |
| $a_{13}^{(r)} + b_{13}^{(r)}$ | $a_{15}^{(r)} + b_{15}^{(r)}$ | $a_{21}^{(r)} + b_{17}^{(r)}$ |
| $a_{14}^{(r)} + b_{14}^{(r)}$ | $a_{16}^{(r)} + b_{16}^{(r)}$ | $a_{22}^{(r)} + b_{18}^{(r)}$ |
| $a_{17}^{(r)} + c_{13}^{(r)}$ | $a_{19}^{(r)} + c_{15}^{(r)}$ | $a_{23}^{(r)} + c_{17}^{(r)}$ |
| $a_{18}^{(r)} + c_{14}^{(r)}$ | $a_{20}^{(r)} + c_{16}^{(r)}$ | $a_{24}^{(r)} + c_{18}^{(r)}$ |
| $b_{19}^{(r)} + c_{19}^{(r)}$ | $b_{21}^{(r)} + c_{21}^{(r)}$ | $b_{23}^{(r)} + c_{23}^{(r)}$ |
| $b_{20}^{(r)} + c_{20}^{(r)}$ | $b_{22}^{(r)} + c_{22}^{(r)}$ | $b_{24}^{(r)} + c_{24}^{(r)}$ |
| $a_{25}^{(r)} + b_{25}^{(r)} + c_{25}^{(r)}$ | $a_{26}^{(r)} + b_{26}^{(r)} + c_{26}^{(r)}$ | $a_{27}^{(r)} + b_{27}^{(r)} + c_{27}^{(r)}$ |

*Round 1:*

$$a_{[1:27]}^{(1)} = \left( W_1^{[1:18]} \mathbf{G}_{[1:18]}^{18 \times 27} + S_{[1:9]}^{(1)} \mathbf{G}_{[19:27]}^{9 \times 27} \right) \mathbf{S}_1^{(1)} \tag{108}$$

$$b_{[1:18]}^{(1)} = \left( W_2^{[19:39]} \mathbf{G}_{[1:21]}^{21 \times 27} + S_{[10:15]}^{(1)} \mathbf{G}_{[22:27]}^{6 \times 27} \right) \mathbf{S}_2^{(1)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{109}$$

$$b_{[19:27]}^{(1)} = \left( W_2^{[19:39]} \mathbf{G}_{[1:21]}^{21 \times 27} + S_{[10:15]}^{(1)} \mathbf{G}_{[22:27]}^{6 \times 27} \right) \mathbf{S}_2^{(1)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{110}$$

$$c_{[1:18]}^{(1)} = \left( W_3^{[40:62]} \mathbf{G}_{[1:23]}^{23 \times 27} + S_{[16:19]}^{(1)} \mathbf{G}_{[24:27]}^{4 \times 27} \right) \mathbf{S}_3^{(1)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{111}$$

$$c_{[19:27]}^{(1)} = \left( W_3^{[40:62]} \mathbf{G}_{[1:23]}^{23 \times 27} + S_{[16:19]}^{(1)} \mathbf{G}_{[24:27]}^{4 \times 27} \right) \mathbf{S}_3^{(1)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{112}$$

*Round 2:*

$$a_{[1:27]}^{(2)} = \left( W_1^{[19:39]} \mathbf{G}_{[1:21]}^{21 \times 27} + S_{[10:15]}^{(2)} \mathbf{G}_{[22:27]}^{6 \times 27} \right) \mathbf{S}_1^{(2)} \tag{113}$$

$$b_{[1:18]}^{(2)} = \left( W_2^{[40:62]} \mathbf{G}_{[1:23]}^{23 \times 27} + S_{[16:19]}^{(2)} \mathbf{G}_{[24:27]}^{4 \times 27} \right) \mathbf{S}_2^{(2)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{114}$$

$$b_{[19:27]}^{(1)} = \left( W_2^{[40:62]} \mathbf{G}_{[1:23]}^{23 \times 27} + S_{[16:19]}^{(2)} \mathbf{G}_{[24:27]}^{4 \times 27} \right) \mathbf{S}_2^{(2)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{115}$$

$$c_{[1:18]}^{(1)} = \left( W_3^{[1:18]} \mathbf{G}_{[1:18]}^{18 \times 27} + S_{[1:9]}^{(2)} \mathbf{G}_{[19:27]}^{9 \times 27} \right) \mathbf{S}_3^{(2)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{116}$$

$$c_{[19:27]}^{(1)} = \left( W_3^{[1:18]} \mathbf{G}_{[1:18]}^{18 \times 27} + S_{[1:9]}^{(2)} \mathbf{G}_{[19:27]}^{9 \times 27} \right) \mathbf{S}_3^{(2)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{117}$$

*Round 3:*

$$a_{[1:27]}^{(3)} = \left( W_1^{[40:62]} \mathbf{G}_{[1:23]}^{23 \times 27} + S_{[16:19]}^{(3)} \mathbf{G}_{[29:27]}^{4 \times 27} \right) \mathbf{S}_1^{(3)} \tag{118}$$

$$b_{[1:18]}^{(3)} = \left( W_1^{[1:18]} \mathbf{G}_{[1:18]}^{18 \times 27} + S_{[1:9]}^{(3)} \mathbf{G}_{[19:27]}^{9 \times 27} \right) \mathbf{S}_2^{(3)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{119}$$

$$b_{[19:27]}^{(3)} = \left( W_1^{[1:18]} \mathbf{G}_{[1:18]}^{18 \times 27} + S_{[1:9]}^{(3)} \mathbf{G}_{[19:27]}^{9 \times 27} \right) \mathbf{S}_2^{(3)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{120}$$

$$c_{[1:18]}^{(3)} = \left( W_3^{[19:39]} \mathbf{G}_{[1:21]}^{21 \times 27} + S_{[10:15]}^{(3)} \mathbf{G}_{[22:27]}^{6 \times 27} \right) \mathbf{S}_3^{(3)}[:, (1:12)] \mathbf{G}_1^{12 \times 18} \tag{121}$$

$$c_{[19:27]}^{(3)} = \left( W_3^{[19:39]} \mathbf{G}_{[1:21]}^{21 \times 27} + S_{[10:15]}^{(3)} \mathbf{G}_{[22:27]}^{6 \times 27} \right) \mathbf{S}_3^{(3)}[:, (13:18)] \mathbf{G}_2^{6 \times 9} \tag{122}$$

**Correctness:** For each round, the user can recover $b_{[13:18]}^{(r)}$ and $c_{[13:18]}^{(r)}$ from $b_{[1:12]}^{(r)}$ and $c_{[1:12]}^{(r)}$. Therefore, the user can cancel the interference and solve $a_{[13:24]}^{(r)}$. Similarly, the user can recover and cancel $b_{[25:27]}^{(r)} + c_{[25:27]}^{(r)}$ and obtain $a_{[25:27]}^{(r)}$, because $b_{[19:27]}^{(r)}$ and $c_{[19:27]}^{(r)}$ are generated from the

same $(9,6)$-MDS code. Hence, the user can solve $a_{[1:27]}^{(r)}$ for all three rounds. For round 1, we have $a_{[1:27]}^{(1)} = \left[W_1^{[1:18]} S_{[1:9]}^{(1)}\right] \mathbf{G} \mathbf{S}_1^{(1)}$. Because $\mathbf{G}$ and $\mathbf{S}_1^{(1)}$ are invertible matrices, the user can solve 18 symbols $W_1^{[1:18]}$. Similarly, the user can solve $W_1^{[19:39]}$ and $W_1^{[40:62]}$ for both round 2 and round 3. Hence, the user obtains all 62 symbols of $W_1$.

**User-privacy:** Any $T = 2$ databases may collude and observe the queries composed of 18 symbols from $a_{[1:27]}^{(r)}$, 12 symbols from both $b_{[1:18]}^{(r)}$ and $c_{[1:18]}^{(r)}$, and 6 symbols from both $b_{[19:27]}^{(r)}$ and $c_{[19:27]}^{(r)}$ for each round. Let $\mathcal{I}_a, \mathcal{I}_{b,12}, \mathcal{I}_{b,6}, \mathcal{I}_{c,12}, \mathcal{I}_{c,6}$ denote the indices of the symbols observed by the colluding databases,

$$
\begin{pmatrix}
a_{\mathcal{I}_a}^{(1)}, & a_{\mathcal{I}_a}^{(2)}, & a_{\mathcal{I}_a}^{(3)} \\
(b_{\mathcal{I}_{b,12}}^{(1)}, b_{\mathcal{I}_{b,6}}^{(1)}), & (b_{\mathcal{I}_{b,12}}^{(2)}, b_{\mathcal{I}_{b,6}}^{(2)}), & (b_{\mathcal{I}_{b,12}}^{(3)}, b_{\mathcal{I}_{b,6}}^{(3)}) \\
(c_{\mathcal{I}_{c,12}}^{(1)}, c_{\mathcal{I}_{c,6}}^{(1)}), & (c_{\mathcal{I}_{c,12}}^{(2)}, c_{\mathcal{I}_{c,6}}^{(2)}), & (c_{\mathcal{I}_{c,12}}^{(3)}, c_{\mathcal{I}_{c,6}}^{(3)})
\end{pmatrix}
\tag{123}
$$

$$
= \begin{pmatrix}
\left[W_1^{[1:18]} S_{[1:9]}^{(1)}\right] \mathbf{GS}_1^{(1)}[:,\mathcal{I}_a], \left[W_1^{[19:39]} S_{[10:15]}^{(2)}\right] \mathbf{GS}_1^{(2)}[:,\mathcal{I}_a], \left[W_1^{[40:62]} S_{[16:19]}^{(3)}\right] \mathbf{GS}_1^{(3)}[:,\mathcal{I}_a] \\
\begin{bmatrix}
\left(\left[W_2^{[19:39]} S_{[10:15]}^{(1)}\right] \mathbf{GS}_2^{(1)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{b,12}], \left[W_2^{[19:39]} S_{[10:15]}^{(1)}\right] \mathbf{GS}_2^{(1)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{b,6}]\right) \\
\left(\left[W_2^{[40:62]} S_{[16:19]}^{(2)}\right] \mathbf{GS}_2^{(2)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{b,12}], \left[W_2^{[40:62]} S_{[16:19]}^{(2)}\right] \mathbf{GS}_2^{(2)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{b,6}]\right) \\
\left(\left[W_2^{[1:18]} S_{[1:9]}^{(3)}\right] \mathbf{GS}_2^{(3)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{b,12}], \left[W_2^{[1:18]} S_{[1:9]}^{(3)}\right] \mathbf{GS}_2^{(3)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{b,6}]\right)
\end{bmatrix}^{\mathrm{T}} \\
\begin{bmatrix}
\left(\left[W_3^{[40:62]} S_{[16:19]}^{(1)}\right] \mathbf{GS}_3^{(1)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{c,12}], \left[W_3^{[40:62]} S_{[16:19]}^{(1)}\right] \mathbf{GS}_3^{(1)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{c,6}]\right) \\
\left(\left[W_3^{[1:18]} S_{[1:9]}^{(2)}\right] \mathbf{GS}_3^{(2)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{c,12}], \left[W_3^{[1:18]} S_{[1:9]}^{(2)}\right] \mathbf{GS}_3^{(2)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{c,6}]\right) \\
\left(\left[W_3^{[19:39]} S_{[10:15]}^{(3)}\right] \mathbf{GS}_3^{(3)}[:,(1:12)] \mathbf{G}_1^{12\times18}[:,\mathcal{I}_{c,12}], \left[W_3^{[19:39]} S_{[10:15]}^{(3)}\right] \mathbf{GS}_3^{(3)}[:,(13:18)] \mathbf{G}_2^{6\times9}[:,\mathcal{I}_{c,6}]\right)
\end{bmatrix}^{\mathrm{T}}
\end{pmatrix}
\tag{124}
$$

$$
\sim \begin{pmatrix}
\left[W_1^{[1:18]} S_{[1:9]}^{(1)}\right] \mathbf{GS}_1^{(1)}[:,(1:18)], \left[W_1^{[19:39]} S_{[10:15]}^{(2)}\right] \mathbf{GS}_1^{(2)}[:,(1:18)], \left[W_1^{[40:62]} S_{[16:19]}^{(3)}\right] \mathbf{GS}_1^{(3)}[:,(1:18)] \\
\begin{bmatrix}
\left(\left[W_2^{[19:39]} S_{[10:15]}^{(1)}\right] \mathbf{GS}_2^{(1)}[:,(1:12)], \left[W_2^{[19:39]} S_{[10:15]}^{(1)}\right] \mathbf{GS}_2^{(1)}[:,(13:18)]\right) \\
\left(\left[W_2^{[40:62]} S_{[16:19]}^{(2)}\right] \mathbf{GS}_2^{(2)}[:,(1:12)], \left[W_2^{[40:62]} S_{[16:19]}^{(2)}\right] \mathbf{GS}_2^{(2)}[:,(13:18)]\right) \\
\left(\left[W_2^{[1:18]} S_{[1:9]}^{(3)}\right] \mathbf{GS}_2^{(3)}[:,(1:12)], \left[W_2^{[1:18]} S_{[1:9]}^{(3)}\right] \mathbf{GS}_2^{(3)}[:,(13:18)]\right)
\end{bmatrix}^{\mathrm{T}} \\
\begin{bmatrix}
\left(\left[W_3^{[40:62]} S_{[16:19]}^{(1)}\right] \mathbf{GS}_3^{(1)}[:,(1:12)], \left[W_3^{[40:62]} S_{[16:19]}^{(1)}\right] \mathbf{GS}_3^{(1)}[:,(13:18)]\right) \\
\left(\left[W_3^{[1:18]} S_{[1:9]}^{(2)}\right] \mathbf{GS}_3^{(2)}[:,(1:12)], \left[W_3^{[1:18]} S_{[1:9]}^{(2)}\right] \mathbf{GS}_3^{(2)}[:,(13:18)]\right) \\
\left(\left[W_3^{[19:39]} S_{[10:15]}^{(3)}\right] \mathbf{GS}_3^{(3)}[:,(1:12)], \left[W_3^{[19:39]} S_{[10:15]}^{(3)}\right] \mathbf{GS}_3^{(3)}[:,(13:18)]\right)
\end{bmatrix}^{\mathrm{T}}
\end{pmatrix}
\tag{125}
$$

$$
= \begin{pmatrix} \left[ W_1^{[1:18]} S_{[1:9]}^{(1)} \right] \mathbf{GS}_1^{(1)}[:, (1:18)], \left[ W_1^{[19:39]} S_{[10:15]}^{(2)} \right] \mathbf{GS}_1^{(2)}[:, (1:18)], \left[ W_1^{[40:62]} S_{[16:19]}^{(3)} \right] \mathbf{GS}_1^{(3)}[:, (1:18)] \\ \left[ W_2^{[19:39]} S_{[10:15]}^{(1)} \right] \mathbf{GS}_2^{(1)}[:, (1:18)], \left[ W_2^{[40:62]} S_{[16:19]}^{(2)} \right] \mathbf{GS}_2^{(2)}[:, (1:18)], \left[ W_2^{[1:18]} S_{[1:9]}^{(3)} \right] \mathbf{GS}_2^{(3)}[:, (1:18)] \\ \left[ W_3^{[40:62]} S_{[16:19]}^{(1)} \right] \mathbf{GS}_3^{(1)}[:, (1:18)], \left[ W_3^{[1:18]} S_{[1:9]}^{(2)} \right] \mathbf{GS}_3^{(2)}[:, (1:18)], \left[ W_3^{[19:39]} S_{[10:15]}^{(3)} \right] \mathbf{GS}_3^{(3)}[:, (1:18)] \end{pmatrix}
$$

$$(126)$$

The user can randomize the three rounds of downloading. Therefore, the symbols requested at the two colluding databases are mapped from the symbols of each file and the $S_i^{(r)}$'s in the same way. Hence, user-privacy is guaranteed.

**System-privacy:** Similar as in the example in Section V-A, the answers from any database is composed by adding linearly independent combinations of $S_{[1:19]}^{(r)}$ for each round. Therefore, the eavesdropper obtains no information regarding the database $W_1, W_2, W_3$ and hence system-privacy is guaranteed.

*F. For arbitrary $N$, $K$, $T$ and $E$ ($E < T$)*

Denote $J = \frac{N^K - T^K}{N - T}$, and suppose each file comprises $L = KN^K - EJ = KN^K - E\frac{N^K - T^K}{N - T}$ symbols from a large enough finite field. The user downloads $K$ rounds, with $NJ$ symbols per round. The database generates $KEJ$ uniformly random symbols, denoted by $S_{[1:EJ]}^{(r)}$ where $r = [1 : K]$.

Divide $[1 : L]$ and $[1 : EJ]$ into $K$ disjoint sets in the following way,

$$
[1 : L] = \underbrace{\mathcal{W}_1}_{\text{size } N^K - EN^{K-1}} \cup \underbrace{\mathcal{W}_2}_{\text{size } N^K - ETN^{K-2}} \cup \ldots \cup \underbrace{\mathcal{W}_{K-1}}_{\text{size } N^K - ET^{K-2}N} \cup \underbrace{\mathcal{W}_K}_{\text{size } N^K - ET^{K-1}} \quad (127)
$$

$$
[1 : EJ] = \underbrace{\mathcal{S}_1}_{\text{size } EN^{K-1}} \cup \underbrace{\mathcal{S}_2}_{\text{size } ETN^{K-2}} \cup \ldots \cup \underbrace{\mathcal{S}_{K-1}}_{\text{size } ET^{K-2}N} \cup \underbrace{\mathcal{S}_K}_{\text{size } ET^{K-1}} \quad (128)
$$

such that $|\mathcal{W}_i| + |\mathcal{S}_i| = N^K$. Therefore, $W_k^{[1:L]} = \{W_k^{\mathcal{W}_1}, \ldots, W_k^{\mathcal{W}_K}\}$ and $S_{[1:EJ]}^{(r)} = \{S_{\mathcal{S}_1}^{(r)}, \ldots, S_{\mathcal{S}_K}^{(r)}\}$.

Let $\{\lambda_1, \ldots, \lambda_{N^K}\}$ be $N^K$ distinct nonzero elements from $\mathbb{F}_q$. Let $\mathbf{G}$ be a $N^K \times N^K$ matrix defined as follows,

$$
\mathbf{G} = \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \lambda_1 & \lambda_2 & \ldots & \lambda_{N^K} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{N^K-1} & \lambda_2^{N^K-1} & \ldots & \lambda_{N^K}^{N^K-1} \end{bmatrix}, \quad (129)
$$

it is direct that $\mathbf{G}$ is an invertible matrix. In the following, we divide $\mathbf{G}$ into $K$ pairs of matrices $\{\mathbf{G}_{\mathcal{W}_i}^{|\mathcal{W}_i| \times N^K}, \mathbf{G}_{\mathcal{S}_i}^{|\mathcal{S}_i| \times N^K}\}$ for $i = [1 : K]$, where $\mathbf{G}_{\mathcal{W}_i}^{|\mathcal{W}_i| \times N^K}$ is composed of the first $|\mathcal{W}_i|$ rows of $\mathbf{G}$ and $\mathbf{G}_{\mathcal{S}_i}^{|\mathcal{S}_i| \times N^K}$ is composed of the last $|\mathcal{S}_i|$ rows of $\mathbf{G}$. It is direct that these $2K$ matrices are generating matrices of MDS codes with corresponding dimensions, and $\mathbf{G} = [\mathbf{G}_{\mathcal{W}_i}^{|\mathcal{W}_i| \times N^K} \; \mathbf{G}_{\mathcal{S}_i}^{|\mathcal{S}_i| \times N^K}]^{\mathrm{T}}$

For each round $r$ and each file index $k$, let $V_k^{(r)}$ be the length-$N^K$ vector defined as follows,

$$V_k^{(r)} = W_k^{\mathcal{W}_{k+r-1 \bmod K}} \mathbf{G}_{\mathcal{W}_{k+r-1 \bmod K}}^{|\mathcal{W}_{k+r-1 \bmod K}| \times N^K} + S_{\mathcal{S}_{k+r-1 \bmod K}}^{(r)} \mathbf{G}_{\mathcal{S}_{k+r-1 \bmod K}}^{|\mathcal{S}_{k+r-1 \bmod K}| \times N^K} \tag{130}$$

$$= \left[ W_k^{\mathcal{W}_{k+r-1 \bmod K}} \; S_{\mathcal{S}_{k+r-1 \bmod K}}^{(r)} \right] \mathbf{G}, \tag{131}$$

therefore, the $K$ index set pairs $(\mathcal{W}_i, \mathcal{S}_i)$ is rotated in all $K$ round for each file index $k \in [1 : K]$. This is to assure user-privacy.

The user privately generates $K^2$ matrices $\mathbf{S}_{[1:K]}^{(1)}, \mathbf{S}_{[1:K]}^{(2)}, \ldots, \mathbf{S}_{[1:K]}^{(K)} \in \mathbb{F}_q^{N^K \times N^K}$ uniformly and independently from all $N^K \times N^K$ invertible matrices over $\mathbb{F}_q$.

Suppose the user wants to retrieve $W_l$. For any undesired file index $k \in [1 : K] \setminus \{l\}$, there are $\Delta = 2^{K-2}$ distinct subsets of $[1 : K]$ which contain $k$ and do not contain $l$, denoted by $\mathcal{K}_1, \mathcal{K}_2, \ldots, \mathcal{K}_\Delta$. For $i \in [1 : \Delta]$, let $\alpha_i = N(N-T)^{|\mathcal{K}_i|-1} T^{K-|\mathcal{K}_i|}$, choose $\Delta$ matrices $\mathbf{G}_1^{\alpha_1 \times \frac{N}{T} \alpha_1}, \ldots, \mathbf{G}_\Delta^{\alpha_\Delta \times \frac{N}{T} \alpha_\Delta}$ be the generating matrices of the MDS codes with corresponding dimensions.

For each round $r$, apply the scheme in [6] for $V_{[1:K]}^{(r)}$ as described in (131). For any undesired file index $k \in [1 : K] \setminus \{l\}$,

$$X_k^{(r)} = \left[ \; x_{\mathcal{K}_1}^{[k],(r)} \; x_{\mathcal{K}_1 \cup \{l\}}^{[k],(r)} \; \Big| \; x_{\mathcal{K}_2}^{[k],(r)} \; x_{\mathcal{K}_2 \cup \{l\}}^{[k],(r)} \; \Big| \; \cdots \; \Big| \; x_{\mathcal{K}_\Delta}^{[k],(r)} \; x_{\mathcal{K}_\Delta \cup \{l\}}^{[k],(r)} \; \right] \tag{132}$$

$$= V_k^{(r)} \mathbf{S}_k^{(r)}[:, (1 : TN^{K-1})] \begin{bmatrix} \mathbf{G}_1^{\alpha_1 \times \frac{N}{T} \alpha_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{G}_2^{\alpha_2 \times \frac{N}{T} \alpha_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{G}_\Delta^{\alpha_\Delta \times \frac{N}{T} \alpha_\Delta} \end{bmatrix}, \tag{133}$$

where the length of $x_{\mathcal{K}_i}^{[k],(r)}$ is $\alpha_i = N(N-T)^{|\mathcal{K}_i|-1} T^{K-|\mathcal{K}_i|}$ and the length of $x_{\mathcal{K}_i \cup \{l\}}^{[k],(r)}$ is $\frac{N-T}{T} \alpha_i$.

For the desired file index $l$, there are $\delta = 2^{K-1}$ distinct subsets of $[1 : K]$ which contain $l$, denoted by $\mathcal{L}_1, \mathcal{L}_2, \ldots, \mathcal{L}_\delta$. Let

$$X_l^{(r)} = \left[ x_{\mathcal{L}_1}^{[l],(r)} \; x_{\mathcal{L}_2}^{[l],(r)} \; \cdots \; x_{\mathcal{L}_\delta}^{[l],(r)} \right] = V_l^{(r)} \mathbf{S}_l^{(r)}, \tag{134}$$

where the length of $x_{\mathcal{L}_i}^{[l],(r)}$ is $N(N-T)^{|\mathcal{L}_i|-1}T^{K-|\mathcal{L}_i|}$.

For each non-empty set $\mathcal{K} \in [1:K]$, the queries associated with $\mathcal{K}$ is generated by

$$\mathcal{Q}_{\mathcal{K}}^{(r)} = \sum_{k \in \mathcal{K}} x_{\mathcal{K}}^{(r)}. \tag{135}$$

For all $K$ rounds $r \in [1:K]$, distribute the queries for each $\mathcal{K}$ evenly among the $N$ databases, and the construction of the queries is completed.

**Decodability, User-privacy, System-privacy, and the Achievable rate**

From [6], for each round, the user can cancel the interference of the undesired files hence obtain $V_l^{(r)}$ for all $K$ rounds. Furthermore, from (131), the user can solve for a different set of symbols $W_l^{\mathcal{W}_i}$ each round, hence the user can obtain all the symbols of the desired file $W_l^{[1:L]} = \{W_l^{\mathcal{W}_1}, \ldots, W_l^{\mathcal{W}_K}\}$.

To see why user-privacy is guaranteed, similarly as in [6], any $T$ colluding servers observe queries comprised of $TN^{K-1}$ symbols of $X_k^{(r)}$ for each round. Denote the index set of $X_k^{(r)}$ observed by the colluding servers by $\mathcal{I}_k$, we have that for all $k \in [1:K]$,

$$X_{\mathcal{I}_k}^{(r)} \sim V_k^{(r)}\mathbf{S}_k^{(r)}[:, (1:TN^{K-1})]. \tag{136}$$

From (131), $V_k^{(r)}$ are constructed from disjoint set of symbols of $W_k$ in an iterative way through the $K$ rounds, and because $\mathbf{S}_k^{(r)}[:, (1:TN^{K-1})]$ are independently and identically distributed, user-privacy is guaranteed since the colluding databases observe symbols constructed from all $W_k$'s through the same random mapping.

System-privacy is guaranteed because from (128) and (131), for each round the $EJ$ queries and answers observed by the eavesdropper is constructed by adding independent linear combinations of $EJ$ independent uniform symbols $S_{[1:EJ]}^{(r)}$. Therefore, the eavesdropper can obtain no information regarding the database $W_{[1:K]}$.

The rate achieved by the scheme is

$$R = \frac{L}{KNJ} = \frac{KN^K - E\frac{N^K-T^K}{N-T}}{KN\frac{N^K-T^K}{N-T}} = \frac{1-\frac{T}{N}}{1-(\frac{T}{N})^K} - \frac{E}{KN} = \underline{R}_{\text{T-EPIR}}. \tag{137}$$

The secrecy rate achieved is

$$\rho = \frac{KEJ}{L} = \frac{KE\frac{N^K-T^K}{N-T}}{KN^K - E\frac{N^K-T^K}{N-T}} = \frac{\frac{E}{N}\left(1-(\frac{T}{N})^K\right)}{1-\frac{T}{N}-\frac{E}{KN}\left(1-(\frac{T}{N})^K\right)}. \tag{138}$$

## REFERENCES

[1] Q. Wang and M. Skoglund, "Secure symmetric private information retrieval from colluding databases with adversaries," *arXiv preprint arXiv:1707.02152*, 2017.

[2] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan, "Private information retrieval," in *IEEE Annual Symposium on Foundations of Computer Science*, 1995, pp. 41–50.

[3] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, "Private information retrieval," *Journal of the ACM (JACM)*, 1998.

[4] Y. Gertner, Y. Ishai, E. Kushilevitz, and T. Malkin, "Protecting data privacy in private information retrieval schemes," in *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, 1998.

[5] H. Sun and S. A. Jafar, "The capacity of private information retrieval," *IEEE Transactions on Information Theory*, 2017.

[6] ——, "The capacity of robust private information retrieval with colluding databases," *arXiv preprint arXiv:1605.00635*, 2016.

[7] ——, "The capacity of symmetric private information retrieval," *arXiv preprint arXiv:1606.08828*, 2016.

[8] K. Banawan and S. Ulukus, "The capacity of private information retrieval from coded databases," *arXiv preprint arXiv:1609.08138*, 2016.

[9] ——, "Multi-message private information retrieval: Capacity results and near-optimal schemes," *arXiv preprint arXiv:1702.01739*, 2017.

[10] ——, "The capacity of private information retrieval from byzantine and colluding databases," *arXiv preprint arXiv:1706.01442*, 2017.

[11] Q. Wang and M. Skoglund, "Symmetric private information retrieval for MDS coded distributed storage," *arXiv preprint arXiv:1610.04530*, 2016.

[12] ——, "Linear symmetric private information retrieval for mds coded distributed storage with colluding servers," *arXiv preprint arXiv:1708.05673*, 2017.

[13] N. B. Shah, K. Rashmi, and K. Ramchandran, "One extra bit of download ensures perfectly private information retrieval," in *Proc. IEEE Int. Symp. Information Theory*, 2014, pp. 856–860.

[14] A. Fazeli, A. Vardy, and E. Yaakobi, "PIR with low storage overhead: coding instead of replication," *arXiv preprint arXiv:1505.06241*, 2015.

[15] T. H. Chan, S.-W. Ho, and H. Yamamoto, "Private information retrieval for coded storage," in *Proc. IEEE Int. Symp. Information Theory*, 2015, pp. 2842–2846.

[16] R. Tajeddine and S. E. Rouayheb, "Private information retrieval from MDS coded data in distributed storage systems," in *Proc. IEEE Int. Symp. Information Theory*, 2016.

[17] R. Freij-Hollanti, O. Gnilke, C. Hollanti, and D. Karpuk, "Private information retrieval from coded databases with colluding servers," *arXiv preprint arXiv:1611.02062*, 2016.

[18] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.