# Multiround Private Information Retrieval: Capacity and Storage Overhead

Hua Sun<sup>10</sup>, Member, IEEE, and Syed Ali Jafar, Fellow, IEEE

Abstract—Private information retrieval (PIR) is the problem of retrieving one message out of K messages from N noncommunicating replicated databases, where each database stores all K messages, in such a way that each database learns no information about which message is being retrieved. The capacity of PIR is the maximum number of bits of desired information per bit of downloaded information among all PIR schemes. The capacity has recently been characterized for PIR as well as several of its variants. In every case it is assumed that all the queries are generated by the user simultaneously. Here we consider multiround PIR, where the queries in each round are allowed to depend on the answers received in previous rounds. We show that the capacity of multiround PIR is the same as the capacity of single-round PIR. The result is generalized to also include T-privacy constraints. Combined with previous results, this shows that there is no capacity advantage from multiround over single-round schemes, non-linear over linear schemes or from  $\epsilon$ -error over zero-error schemes. However, we show through an example that there is an advantage in terms of storage overhead. We provide an example of a multiround, non-linear,  $\epsilon$ -error PIR scheme that requires a strictly smaller storage overhead than the best possible with single-round, linear, zero-error PIR schemes.

*Index Terms*—Private information retrieval, multiple rounds, capacity, storage overhead.

#### I. INTRODUCTION

**P**RIVATE information retrieval (PIR) [1], [2] is one of the canonical problems in theoretical computer science and cryptography. The PIR setting involves K messages that are assumed to be independent, N distributed databases that are replicated (each database stores all K messages) and noncolluding (the databases do not communicate with each other), and a user who desires one of the K messages. A PIR scheme is any mechanism by which a user may retrieve his desired message from the databases privately, i.e., without revealing any information about which message is being retrieved, to any

Manuscript received February 2, 2017; revised July 31, 2017; accepted December 16, 2017. Date of publication January 4, 2018; date of current version July 12, 2018. This work was supported in part by NSF under Grant CCF-1617504 and Grant CNS-1731384, in part by ONR under Grant N00014-16-1-2629 and Grant N00014-18-1-2057, and in part by ARL under Grant W911NF-16-1-0215.

H. Sun is with the Department of Electrical Engineering, University of North Texas, Denton, TX 76203 USA (e-mail: hua.sun@unt.edu).

S. A. Jafar is with the Center for Pervasive Communications and Computing, Department of Electrical Engineering and Computer Science, University of California at Irvine, Irvine, CA 92697 USA (e-mail: syed@uci.edu).

Copyright © 2017 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from

the IEEE by sending a request to pubs-permissions@ieee.org. Communicated by A. Khisti, Associate Editor for Shannon Theory.

Digital Object Identifier 10.1109/TIT.2018.2789426

individual database. An information theoretic formulation of PIR guarantees the user's privacy even if the databases are computationally unbounded.<sup>1</sup> The "rate" of a PIR scheme is defined as the ratio of the number of bits of desired information to the total number of bits downloaded by the user from all the databases. The supremum of achievable rates is defined to be the capacity of PIR. For *K* messages and *N* databases, the capacity of PIR was characterized recently in [6] as

$$C = \left(1 + 1/N + 1/N^2 + \dots + 1/N^{K-1}\right)^{-1}$$
(1)

The capacity has also been determined for various constrained forms of PIR such as LPIR [7] – where message *lengths* can be arbitrary, TPIR [8] – where any set of up to *T* databases may collude, RPIR [8] – where *robustness* is required against unresponsive databases, SPIR [9] – which extends the privacy constraint *symmetrically* to protect both the user and the databases, MDS-PIR [10] and MDS-SPIR [11] – variants of PIR and SPIR, respectively, where each message is separately MDS coded.<sup>2</sup>

A common theme in these results is that there is no capacity advantage of non-linear schemes over linear schemes, or of  $\epsilon$ -error schemes over zero-error schemes. This is a matter of some curiosity because the necessity of non-linear coding schemes has often been a key obstacle in network coding capacity problems [13]–[16], and the capacity benefit of  $\epsilon$ error schemes over zero-error schemes for network coding problems in general [17] remains one of the key unresolved mysteries — with direct connections to the edge-removal question [18] and the existence of strong converses [19] in network information theory. Motivated by this curiosity, in this work we explore another important variant of PIR – *multiround* PIR (MPIR). Our contributions are summarized next.

The classical PIR setting assumes that all the queries are simultaneously generated by the user. This assumption is also made in [6]. However, such a constraint is not essential to PIR. What if this constraint is relaxed, i.e., multiple rounds of queries and answers are allowed, such that the queries

0018-9448 © 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

<sup>&</sup>lt;sup>1</sup>There is also a widely studied cryptographic formulation of PIR, where the user's privacy is guaranteed only against computationally bounded databases [3]–[5].

<sup>&</sup>lt;sup>2</sup>As a caveat, we note that separate MDS coding of each message is a restrictive assumption. Consider the setting with K = 2 messages, N = 3 databases and the storage size of each database is equal to the size of one message. If separate MDS codes are employed for each message, then the maximum rate (capacity) is equal to 3/5 [10]. However, [12, Example 2] shows that rate 2/3 (> 3/5) is achievable with a storage code that jointly encodes both messages.

in each round of communication are generated by the user with the knowledge of the answers from all previous rounds? The resulting variant of the PIR problem is the *multiround* PIR (MPIR) problem (also known as interactive PIR [20], [21]). Multiround PIR has been noted as an intriguing possibility in several prior works [2], [20], [21]. However, it is not known whether there is any benefit of MPIR over single-round PIR. Answering this question from a capacity perspective is the first contribution of this work. Specifically, we show that the capacity of MPIR is the same as the capacity of PIR, i.e., both are given by (1). Combined with previous results, this shows that there is no capacity advantage from multiround over single-round schemes, non-linear over linear schemes or from  $\epsilon$ -error over zero-error schemes. Furthermore, we show that this is true even with *T*-privacy constraints.

To complement the capacity analysis, we consider another metric of interest - storage overhead. Classical PIR assumes replicated databases, i.e., each database stores all the messages. For larger datasets, replication schemes incur substantial storage costs. Coding has been shown to be an effective way to reduce the storage costs in distributed data storage systems. Applications of coding to reduce the storage overhead for PIR have attracted attention recently [10]–[12], [22]–[28]. In this context, our main contribution is an example (N = 2)databases, K = 2 messages) of a multiround, non-linear,  $\epsilon$ -error PIR scheme that achieves a strictly smaller storage overhead than the best possible with a single-round, linear, zero-error scheme. The simplicity of the scheme and the N = K = 2 setting makes it an attractive point of reference for future work toward understanding the role of linear versus non-linear schemes, zero-error versus  $\epsilon$ -error capacity, and single-round versus multiround communications. Interestingly, the scheme reveals that coded storage is useful not only for reducing the storage overhead, but also it has a surprising benefit of enhancing the privacy of PIR.

Notation: For  $n_1, n_2 \in \mathbb{Z}$ ,  $n_1 \leq n_2$ , define the notation  $[n_1 : n_2]$  as the set  $\{n_1, n_1 + 1, \dots, n_2\}$ ,  $A(n_1 : n_2)$  as the vector  $(A(n_1), A(n_1 + 1), \dots, A(n_2))$  and  $A_{n_1:n_2}$  as the vector  $(A_{n_1}, A_{n_1+1}, \dots, A_{n_2})$ . In this paper, we follow the convention that for queries and answers, sub-scripts denote the database index, super-scripts denote the message index and parentheses denote the communication round index. When  $n_1 > n_2$ ,  $[n_1 : n_2]$  is a null set and  $A(n_1 : n_2)$ ,  $A_{n_1:n_2}$  are null vectors. For an index set  $T = \{i_1, i_2, \dots, i_n\}$  such that  $i_1 < i_2 < \dots < i_n$ , the notation  $A_T$  represents the vector  $(A_{i_1}, A_{i_2}, \dots, A_{i_n})$ . The notation  $X \sim Y$  is used to indicate that X and Y are identically distributed.

## **II. PROBLEM STATEMENT**

Let us start with a general problem statement that can then be specialized to various settings of interest. Consider K independent messages  $W_1, \dots, W_K$ , each comprised of Li.i.d. uniform bits.

$$H(W_1, \cdots, W_K) = H(W_1) + \cdots + H(W_K),$$
 (2)

$$H(W_1) = \dots = H(W_K) = L. \tag{3}$$

There are N databases. Let  $S_n$  denote the random variable that represents the information stored at the  $n^{th}$  database.

$$H(S_n|W_1, W_2, \cdots, W_K) = 0, \quad \forall n \in [1:N].$$
 (4)

Define the storage overhead  $\alpha$  as the ratio of the total amount of storage used by all databases to the total amount of data.<sup>3</sup>

$$\alpha \stackrel{\triangle}{=} \frac{\sum_{n=1}^{N} H(S_n)}{KL}.$$
(5)

For replication based schemes, each database stores all K messages, so  $S_n = (W_1, W_2, \dots, W_K), H(S_n) = KL$ , and the storage overhead,  $\alpha = N$ .

A user privately generates  $\theta$  uniformly from [1 : K] and wishes to retrieve  $W_{\theta}$  while keeping  $\theta$  a secret from each database.

Prior works on capacity of PIR and its variants make certain (implicitly justified) assumptions of deterministic behavior, e.g., that the answers provided by the databases are deterministic functions of queries and messages. Here we will follow, instead, an explicit formulation. We allow randomness in the strategies followed by the user and the databases. This is accomplished by representing the actions of the user and the databases as functions of random variables. Let us use  $\mathbb{F}$  to denote a random variable privately generated by the user, whose realization is not available to the databases (the distribution of  $\mathbb{F}$  could be made public to all databases). Similarly, G is a random variable that determines the random strategies followed by the databases, and whose realizations are assumed to be known to all the databases<sup>4</sup> and the user without loss of generality.  $\mathbb{F}$  and  $\mathbb{G}$  take values over the set of all deterministic strategies that the user or the databases can follow, respectively, associating each strategy with a certain probability.  $\mathbb{F}$  and  $\mathbb{G}$  are generated offline, i.e., before the realizations of the messages or the desired message index are known. Since these random variables are generated a-priori, we must have

$$H(\theta, \mathbb{F}, \mathbb{G}, W_1, \cdots, W_K)$$
  
=  $H(\theta) + H(\mathbb{F}) + H(\mathbb{G}) + H(W_1) + \cdots + H(W_K).$  (6)

The multiround PIR scheme proceeds as follows. Suppose  $\theta = k$ . In order to retrieve  $W_k, k \in [1 : K]$  privately, the user communicates with the databases over  $\Gamma$  rounds. In the first round, the user privately generates N random queries,  $Q_1^{[k]}(1), Q_2^{[k]}(1), \cdots, Q_N^{[k]}(1)$ .

$$H(Q_1^{[k]}(1), Q_2^{[k]}(1), \cdots, Q_N^{[k]}(1)|\mathbb{F}) = 0, \quad \forall k \in [1:K].$$
(7)

The user sends query  $Q_n^{[k]}(1)$  to the  $n^{th}$  database,  $\forall n \in [1 : N]$ . Upon receiving  $Q_n^{[k]}(1)$ , the  $n^{th}$  database generates an answering string  $A_n^{[k]}(1)$ . Without loss of generality, we assume

<sup>&</sup>lt;sup>3</sup>Perfect compression may not be possible for arbitrary *L*, especially when *L* is small. However, since in this work we consider the Shannon theoretic formulation where the message size is large, i.e.,  $L \rightarrow \infty$ , we have defined storage overhead using the entropy of *S<sub>n</sub>* which is achievable in this regime.

<sup>&</sup>lt;sup>4</sup>One might wonder if we could allow each database to have its own random strategy (determined by a random variable  $\mathbb{G}_i$ ). We note that in this case, we can define  $\mathbb{G} = (\mathbb{G}_1, \cdots, \mathbb{G}_N)$  such that the constraints in this paper all continue to hold (i.e., (8), (10), (11), (13)) and the converse in Theorem 1 still applies. Thus, there is no potential gain of localized  $\mathbb{G}_i$ .

that the answering string is a function of  $Q_n^{[k]}(1)$ , the stored information  $S_n$ , and the random variable  $\mathbb{G}$ .

$$H(A_n^{[k]}(1)|Q_n^{[k]}(1), S_n, \mathbb{G}) = 0.$$
(8)

Each database returns to the user its answer  $A_n^{[k]}(1)$ .

Proceeding similarly,<sup>5</sup> over the  $\gamma^{th}$  round,  $\gamma \in [2 : \Gamma]$ , the user generates N queries  $Q_1^{[k]}(\gamma), \dots, Q_N^{[k]}(\gamma)$ , which are functions of previous queries and answers and  $\mathbb{F}$ ,

$$H(Q_{1:N}^{[k]}(\gamma)|Q_{1:N}^{[k]}(1:\gamma-1), A_{1:N}^{[k]}(1:\gamma-1), \mathbb{F}) = 0.$$
(9)

The user sends query  $Q_n^{[k]}(\gamma)$  to the  $n^{th}$  database, which generates an answer  $A_n^{[k]}(\gamma)$  and returns  $A_n^{[k]}(\gamma)$  to the user. The answer is a function of all queries received so far, the stored information  $S_n$ , and  $\mathbb{G}$ ,

$$H(A_n^{[k]}(\gamma)|Q_n^{[k]}(1:\gamma), S_n, \mathbb{G}) = 0.$$
(10)

At the end of  $\Gamma$  rounds, from all the information that is now available to the user  $(A_{1:N}^{[k]}(1 : \Gamma), Q_{1:N}^{[k]}(1 : \Gamma), \mathbb{F})$ , the user decodes the desired message  $W_k$  according to a decoding rule that is specified by the PIR scheme. Let  $P_e$  denote the probability of error achieved with the specified decoding rule.

To protect the user's privacy, the K possible values of the desired message index should be indistinguishable from the perspective of any subset  $\mathcal{T} \subset [1 : N]$  of at most T colluding databases, i.e., the following privacy constraint must be satisfied.

$$[T - Privacy] (\mathcal{Q}_{\mathcal{T}}^{[k]}(1:\Gamma), \mathcal{A}_{\mathcal{T}}^{[k]}(1:\Gamma), \mathbb{G}, S_{\mathcal{T}}) \sim (\mathcal{Q}_{\mathcal{T}}^{[k']}(1:\Gamma), \mathcal{A}_{\mathcal{T}}^{[k']}(1:\Gamma), \mathbb{G}, S_{\mathcal{T}}) \forall k, k' \in [1:K], \quad \forall \mathcal{T} \subset [1:N], |\mathcal{T}| = T.$$

$$(11)$$

The PIR rate characterizes how many bits of desired information are retrieved per downloaded bit and is defined as follows:<sup>6</sup>

$$R = \frac{L}{D} \tag{12}$$

where *D* is the expected value<sup>7</sup> of the total number of bits downloaded by the user from all the databases over all  $\Gamma$  rounds.

<sup>5</sup>One might wonder if the setting can be further generalized by allowing sequential queries, i.e., allowing the query to each database to depend not only on the answers received from previous rounds, but also on the answers received from other databases queried previously within the same round. We note that sequential queries are already contained in our multiround framework, e.g., by querying only one database in each round (sending null queries to the remaining databases).

<sup>6</sup>In this work, the metric we consider is rate (download cost), while the upload cost is ignored. We note that in the single round setting [6], the upload cost is negligible in the large message size regime because the same query can be reused multiple times. However, in the multi-round setting as considered in this work, the queries depend on previous answers so that the same query may not be reused. Therefore, as a caveat we note that for multi-round schemes the upload cost may not be negligible.

<sup>7</sup>Alternatively, *D* may be defined as the maximum download needed by the PIR scheme which (similar to choosing zero-error instead of  $\epsilon$ -error) weakens the converse and strengthens the achievability arguments in general. The capacity characterizations in this work, as well as previous works in [6], [8], and [9] hold under either definition. This is because in every case, the upper bounds allow average download *D*, while the achievability only requires maximum download *D*.

A rate *R* is said to be  $\epsilon$ -error achievable if there exists a sequence of PIR schemes, indexed by *L*, each of rate greater than or equal to *R*, for which  $P_e \to 0$  as  $L \to \infty$ .<sup>8</sup> Note that for such a sequence of PIR schemes, from Fano's inequality we must have  $\forall k \in [1:K]$ :

$$o(L) = H(W_k | A_{1:N}^{[k]}(1:\Gamma), Q_{1:N}^{[k]}(1:\Gamma), \mathbb{F}, \mathbb{G})$$

$$\stackrel{(7)(9)}{=} H(W_k | A_{1:N}^{[k]}(1:\Gamma), \mathbb{F}, \mathbb{G}), \qquad (13)$$

where any function of L, say f(L) is said to be o(L) if  $\lim_{L\to\infty} f(L)/L = 0$ . The supremum of  $\epsilon$ -error achievable rates is called the  $\epsilon$ -error capacity  $C_{\epsilon}$ .

A rate *R* is said to be zero-error achievable if there exists (for some *L*) a PIR scheme of rate greater than or equal to *R* for which  $P_e = 0$ . The supremum of zero-error achievable rates is called the zero-error capacity  $C_o$ . From the definitions, it is evident that

$$C_o \le C_\epsilon. \tag{14}$$

### III. RESULTS

There are two main contributions in this work, summarized in the following sections.

#### A. Capacity Perspective

We first consider the capacity benefits of multiple rounds of communication in the classical setting where each database stores all messages, i.e., storage is unconstrained. We present our result in the general context of multiround PIR with T-privacy constraints (MTPIR). The MTPIR setting is obtained from the general problem statement by relaxing the storage overhead constraints, i.e.,

$$S_n = (W_1, W_2, \cdots, W_K), \forall n \in [1:N]$$
  
$$\alpha = N$$

i.e., each database stores all the messages (replication). The following theorem presents the main result.

Theorem 1: The capacity of MTPIR

$$C_o = C_{\epsilon} = \left(1 + T/N + T^2/N^2 + \dots + T^{K-1}/N^{K-1}\right)^{-1}.$$

The converse proof of Theorem 1 is presented in Section IV. Achievability follows directly from [8]. The following observations place the result in perspective.

- The capacity of MTPIR matches the capacity of TPIR found in [8], i.e., multiple rounds do not increase capacity.
- 2) Setting T = 1 gives us the capacity of multiround PIR (MPIR) without *T*-privacy constraints. The capacity of MPIR matches the capacity of PIR found in [6], i.e., multiple rounds do not increase capacity.
- 3) Since the achievability proofs in [6] and [8] only require linear and zero-error schemes, there is no capacity

<sup>8</sup>Equivalently, for any  $\epsilon > 0$ , there exists a finite  $L_{\epsilon}$  such that  $P_e < \epsilon$  for all  $L > L_{\epsilon}$ .

benefit of multiple rounds over single-round schemes, non-linear over linear schemes, or  $\epsilon$ -error over zero-error schemes.

- For all N, K, T, Γ the converse proof of Theorem 1 generalizes the converse proofs of [6] and [8]. Remarkably, it requires only Shannon information inequalities, i.e., sub-modularity of entropy.
- 5) Since the capacity metric focuses only on the download cost, it does not penalize multi-round PIR schemes for their potentially non-negligible upload cost relative to single-round PIR schemes. Theorem 1 shows that even with the ability to have unlimited uploads for free (which particularly favors multi-round schemes) multi-round PIR offers no capacity advantage over single-round PIR.

#### B. Storage Overhead Perspective

As summarized above, our first result shows that in a broad sense - with or without colluding databases - there is no capacity benefit of multiple rounds over single-round communication,  $\epsilon$ -error over zero-error schemes or non-linear over linear schemes for PIR. This pessimistic finding may lead one to believe that there is little reason to further explore interactive communication, non-linear schemes or  $\epsilon$ -error schemes for PIR. As our main contribution in this section, we offer an optimistic counterpoint by looking at the PIR problem from the perspective of storage overhead instead of capacity. The counterpoint is made through a counterexample. The counterexample is quite remarkable in itself as it shows from a storage overhead perspective not only the advantage of a multiround PIR scheme over all single-round PIR schemes, but also of a non-linear PIR scheme over all linear PIR schemes, and an  $\epsilon$ -error scheme over all zero-error schemes.

For a counterexample the simplest setting is typically the most interesting. Therefore, in this section we will only consider the simplest non-trivial setting, with K = 2 messages, N = 2 databases, and T = 1, i.e., no collusion among databases. Recall that for this setting the capacity is C = 2/3. For our counterexample we explore the minimum storage overhead that is needed to achieve the rate 2/3.

Theorem 2: For K = 2, N = 2, T = 1, and for rate 2/3,

1) there exists a multiround, non-linear and  $\epsilon$ -error PIR scheme with storage overhead

$$\alpha = 3/4 + 3/8 \log_2 3$$

which is less than 3/2.

2) the storage overhead of any single-round, linear and zero-error PIR scheme is

$$\alpha \geq 3/2$$

The achievability arguments, including the multiround, nonlinear and  $\epsilon$ -error PIR scheme that proves the first part of Theorem 2 are presented in this section. The proof of the second claim notably utilizes Ingleton's inequality, which goes beyond submodularity, and is presented in Section V. 1) A Multiround, Non-Linear, and  $\epsilon$ -Error PIR Scheme for K = 2, N = 2, T = 1: Define  $w_1, w_2$  as two independent uniform binary random variables. Further, define

$$x_1 = w_1 \wedge w_2 \tag{15}$$

$$x_2 = (\neg w_1) \land \ (\neg w_2) \tag{16}$$

$$y_1 = w_1 \land \ (\neg w_2) \tag{17}$$

$$y_2 = (\neg w_1) \land w_2 \tag{18}$$

where  $\wedge$  and  $\neg$  are the logical AND and NOT operators. Note the following,

$$x_1 = 1 \Rightarrow (w_1, w_2) = (1, 1) \tag{19}$$

$$x_2 = 1 \Rightarrow (w_1, w_2) = (0, 0) \tag{20}$$

$$x_1 = 0 \Rightarrow (w_1, w_2) = (y_1, y_2)$$
 (21)

$$x_2 = 0 \Rightarrow (w_1, w_2) = (\neg y_2, \neg y_1)$$
 (22)

For ease of exposition, consider first the case where each message is only one bit long. In this case, the messages  $W_1, W_2$ , directly correspond to  $w_1, w_2$ , respectively. Denote the first database as DB1 and the second database as DB2. Regardless of whether the user desires  $W_1$  or  $W_2$ , he flips a private fair coin, and requests the value of either  $x_1$  (head) or  $x_2$  (tail) from DB1. If the answer is 1, then according to (19) and (20) the user knows the values of both  $w_1, w_2$  and no further information is requested from DB2. If the answer is 0, then the user proceeds as follows.

- If  $x_1 = 0$  and  $W_1$  is desired, ask DB2 for the value of  $y_1$ . Retrieve  $w_1 = y_1$ .
- If  $x_1 = 0$  and  $W_2$  is desired, ask DB2 for the value of  $y_2$ . Retrieve  $w_2 = y_2$ .
- If x<sub>2</sub> = 0 and W<sub>1</sub> is desired, ask DB2 for the value of y<sub>2</sub>. Retrieve w<sub>1</sub> = ¬y<sub>2</sub>.
- If  $x_2 = 0$  and  $W_2$  is desired, ask DB2 for the value of  $y_1$ . Retrieve  $w_2 = \neg y_1$ .

Note that in order to answer the user's queries, DB1 only needs to store  $(x_1, x_2)$ , and DB2 only needs to store  $(y_1, y_2)$ . This observation is the key to not only the reduced storage overhead, but also the enhanced privacy of this scheme.

Further, in preparation for the proofs that follow, let us define another binary random variable u, which takes the value u = 0 if no response is needed from DB2, and the value u = 1 otherwise. Note that u = 0 implies that  $(y_1, y_2) = (0, 0)$ . On the other hand, if u = 1, then  $(y_1, y_2)$  takes the values (0, 0), (1, 0), (0, 1), each with probability 1/3. Therefore,

 $H(y_1, y_2|u)$ 

$$= 1/4 \times H(y_1, y_2|u=0) + 3/4 \times H(y_1, y_2|u=1) \quad (23)$$

$$= 1/4 \times 0 + 3/4 \times H(1/3, 1/3, 1/3) = 3/4 \log_2 3 \quad (24)$$

The correctness of the scheme is obvious from (19)-(22). Let us verify that the scheme is private. Start with DB1. The query to DB1 is equally likely to be  $x_1$  or  $x_2$ , regardless of the desired message index and the message realizations. Therefore, DB1 learns nothing about which message is retrieved. Next consider DB2. Let us prove that  $(Q_2^{[1]}, y_1, y_2) \sim$ 

$$\begin{array}{c|c} (\mathcal{Q}_{2}^{[2]}, y_{1}, y_{2}). \\ \hline & (\theta = 1) \\ \hline (\mathcal{Q}_{2}^{[1]}, y_{1}, y_{2}) & \underline{\text{Prob.}} \\ \hline & (\mathcal{Q}_{2}^{[2]}, y_{1}, y_{2}) & \underline{\text{Prob.}} \\ \hline$$

where the double quote notation around a random variable represents the query about its realization. The computation of the joint distribution values is straightforward. We present the derivation here for one case. All other cases follow similarly. From the law of total probability, we have

$$Pr\left(\left(Q_{2}^{[1]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)\right)$$

$$= Pr\left(\left(Q_{2}^{[1]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)| (Q_{1}^{[1]}, w_{1}, w_{2})$$

$$= ("x_{1}", 0, 1)\right)$$

$$\times Pr\left(\left(Q_{1}^{[1]}, w_{1}, w_{2}\right) = ("x_{1}", 0, 1)\right)$$

$$+ Pr\left(\left(Q_{2}^{[1]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)| (Q_{1}^{[1]}, w_{1}, w_{2})$$

$$= ("x_{2}", 0, 1)\right)$$

$$\times Pr\left(\left(Q_{1}^{[1]}, w_{1}, w_{2}\right) = ("x_{2}", 0, 1)\right)$$

$$(25)$$

$$= 1 \times 1/8 + 0 \times 1/8 = 1/8$$

$$(26)$$

Similarly,

$$Pr\left(\left(\mathcal{Q}_{2}^{[2]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)\right)$$

$$= Pr\left(\left(\mathcal{Q}_{2}^{[2]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)\right) \left(\mathcal{Q}_{1}^{[2]}, w_{1}, w_{2}\right)$$

$$= ("x_{1}", 0, 1)\right)$$

$$\times Pr\left(\left(\mathcal{Q}_{1}^{[2]}, w_{1}, w_{2}\right) = ("x_{1}", 0, 1)\right)$$

$$+ Pr\left(\left(\mathcal{Q}_{2}^{[2]}, y_{1}, y_{2}\right) = ("y_{1}", 0, 1)\right) \left(\mathcal{Q}_{1}^{[2]}, w_{1}, w_{2}\right)$$

$$= ("x_{2}", 0, 1)\right)$$

$$\times Pr\left(\left(\mathcal{Q}_{1}^{[2]}, w_{1}, w_{2}\right) = ("x_{2}", 0, 1)\right) (27)$$

$$= 0 \times 1/8 + 1 \times 1/8 = 1/8$$
 (28)

Thus,  $\Pr\left((Q_2^{[1]}, y_1, y_2) = ("y_1", 0, 1)\right) = \Pr\left((Q_2^{[1]}, y_1, y_2) = ("y_1", 0, 1)\right)$ . All other cases are verified similarly. Then, since the distribution of  $(Q_2^{[\theta]}, y_1, y_2)$  does not depend on  $\theta$ , and the answers are only deterministic functions of the query and the stored information, it follows that the scheme is private.

Next consider the *L* length extension of this PIR scheme, where each desired bit is retrieved independently as described above. Under the *L* length extension,  $W_1, W_2, X_1, X_2, Y_1, Y_2, U$  are sequences of length *L*, such that the sequence of tuples  $[(W_1(l), W_2(l), X_1(l), X_2(l), Y_1(l), Y_2(l), U(l))]_{l=1}^L$  is i.i.d.  $\sim (w_1, w_2, x_1, x_2, y_1, y_2, u)$ .

Since the extension is obtained by repeated independent applications of the PIR scheme described above for retrieving each message bit, it follows trivially that the extended PIR scheme is also correct and private. The purpose for the *L* length extension, with  $L \rightarrow \infty$ , is to invoke fundamental limits of data compression which optimize both the data rates and the storage overhead as explained next.

Let us show that the rate 2/3 is achieved asymptotically as  $L \to \infty$ . We take advantage of the fact that the answers from the databases are not uniformly distributed, and therefore the sequence of answers from each database is compressible (i.e., each database codes over the sequence of answers and returns the compressed answer to the user). With optimal compression, the user downloads H(1/4, 3/4) bits per desired message bit from DB1. This is because, for each retrieved bit, the answer from DB1 takes the value 1 with probability 1/4 and 0 with probability 3/4. From DB2, we download  $1/4 \times 0 + 3/4 \times H(1/3, 2/3) = 3/4 H(1/3, 2/3)$  bits per desired message bit, because with probability 1/4 (when the answer from DB1 is 1), no response is requested from DB2 and otherwise within the remaining space of probability measure 3/4 (when the answer from DB1 is 0), the answer from DB2 is 1 with conditional probability 1/3 and 0 with conditional probability 2/3. Therefore the total download is H(1/4, 3/4) + 3/4 H(1/3, 2/3) = 3/2 bits per desired message bit and the rate achieved is 2/3.

Next let us determine the storage requirements of this scheme. DB1 needs  $(X_1, X_2)$  to answer the user's queries, so with optimal compression, it needs to store  $H(x_1, x_2) = H(1/4, 1/4, 1/2) = 3/2$  bits per desired message bit. One might naively imagine that the same storage requirement also applies to DB2, because DB2 similarly needs the values  $(Y_1, Y_2)$  to answer the user's queries. However, this is not true, because the query sent to DB2 already contains some information about the message realizations,<sup>9</sup> and this *side-information* allows DB2 to reduce its storage requirement by taking advantage of Slepian Wolf coding [29], [30] (distributed compression with decoder side information).

The key is to realize that DB2 does not need to know  $(Y_1, Y_2)$  until after it receives the query from the user. The query from the user includes U as side information. Therefore, using Slepian Wolf coding, DB2 is able to optimally compress the i.i.d. sequence  $(Y_1, Y_2)$  to the conditional entropy  $H(y_1, y_2|u)$  bits per desired message bit and still decode the  $(Y_1, Y_2)$  sequence when it is needed, i.e., after the query is provided by the user. Thus, the total storage required by this PIR scheme is  $3/2 + 3/4 \log_2 3$  bits per bit of desired message. Since there are two messages, the storage overhead is  $3/4 + 3/8 \log_2 3$ .

The following observations are useful to place the new PIR scheme in perspective.

1) The optimal compression guarantees (Slepian Wolf coding) are only available in the  $\epsilon$ -error sense. Therefore, this PIR scheme is essentially an  $\epsilon$ -error scheme.

<sup>&</sup>lt;sup>9</sup>Note that the query sent to DB2 is independent of the desired message index but not the message realizations.

- 2) The multiround scheme is in fact a sequential PIR scheme that utilizes only one round of queries for each database (two rounds total since there are two databases).
- 3) The scheme is essentially non-linear because, e.g., the logical AND operator is non-linear.
- 4) Since the multiround, non-linear and  $\epsilon$ -error aspects are all essential for *this* scheme to create an advantage in terms of storage overhead, it is an intriguing question whether all three aspects are necessary in *general* or if it is possible to achieve storage overhead less than 3/2 through another scheme while sacrificing at least one of the three aspects.
- 5) A key insight from this PIR scheme is the surprising privacy benefit of storage overhead optimization. By not storing all the information at each database, and by optimally compressing the stored information, not only do we reduce the storage overhead, but also we enable stronger privacy guarantees than would hold otherwise. Note that if each database stores all the information (both  $W_1$  and  $W_2$ ), then the scheme is not private. To see this, suppose  $(w_1, w_2) = (1, 1)$ . This would be known to DB2 because it stores both messages. Under this circumstance, DB2 knows that if the user asks for  $y_2$ , then his desired message must be  $W_1$  and if the user asks for  $y_1$  then his desired message must be  $W_2$ . Thus, storing all the information at each database would result in loss of privacy. As another example, we note that if  $w_1$  and/or  $w_2$  are not uniformly distributed then again the PIR scheme would lose privacy. To see this, suppose  $Pr(w_1 = 1) = Pr(w_2 = 1) > 1/2$ . Then DB2 is more likely to be asked for  $y_1$  if the desired message is  $W_2$  than if the desired message is  $W_1$ . On the other hand, note that optimal data compression is a pre-requisite in any case for the optimization of rate and storage overhead.10
- 6) Let us consider momentarily the restricted message size setting, where each message is only L = 1 bit long. Then it is easy to see that any single-round scheme (all queries generated simultaneously) must download at least 2 bits on average, but our multiround scheme requires an expected download of only 1 + 3/4 = 7/4 bits. Thus, when storage is not constrained, even though the download advantage of multiround PIR disappears under unconstrained message lengths, for constrained message lengths there are benefits of multiround PIR.
- 7) A key limitation of our PIR scheme is its upload cost. We need to send queries for each message bit so the upload cost scales linearly with the message size *L*. This observation points to an interesting tradeoff between the storage overhead and upload cost. In particular the possibility of storage overhead improvements subject to negligible upload cost is intriguing.

2) A Single-Round, Linear, and Zero-Error Scheme for K = 2, N = 2, T = 1: For comparison, the corresponding scheme from [6] which also achieves rate 2/3 is reproduced below. This will be shown to be the optimal single-round, linear, zero-error scheme for storage overhead in Section V. Denote the messages, each comprised of 4 bits, as  $W_1 = (a_1, a_2, a_3, a_4), W_2 = (b_1, b_2, b_3, b_4)$ . The downloaded information from each database is shown at the top of the next page.

The scheme achieves rate 2/3 and is linear, single-round, and zero-error. A total of 6 bits are stored at each database

$$S_1 = (a_1, a_3, b_1, b_3, a_2 + b_2, a_4 + b_4)$$
(29)

$$S_2 = (a_2, a_4, b_2, b_4, a_3 + b_1, a_1 + b_3)$$
(30)

Thus, the storage overhead is 3/2.

## IV. PROOF OF THEOREM 1

We first present two useful lemmas. Note that in the proofs, the relevant equations needed to justify each step are specified by the equation numbers set on top of the (in)equality symbols. *Lemma 1:* For all  $k \in [2: K]$ ,

$$I(W_{k:K}; Q_{1:N}^{[k-1]}(1:\Gamma), A_{1:N}^{[k-1]}(1:\Gamma), \mathbb{F}|W_{1:k-1}, \mathbb{G}) \geq \frac{T}{N} I(W_{k+1:K}; Q_{1:N}^{[k]}(1:\Gamma), A_{1:N}^{[k]}(1:\Gamma), \mathbb{F}|W_{1:k}, \mathbb{G}) + \frac{LT}{N} (1 - \frac{o(L)}{L}).$$
(31)

*Proof:* The proof is shown at the top of the next page.

Here, (32) follows from the non-negativity of mutual information. In (33), we have used the privacy constraint (11) and in this storage unconstrained setting, the stored information  $S_T$  in (11) is  $W_{1:K}$ . (34) is due to the chain rule and the fact that mutual information is non-negative. In (37), we use Han's inequality [30, Th. 17.6.1] with conditioning on common random variables.

*Remark:* Intuitively, Lemma 1 recursively relates the interference (information about other messages that is contained in the answers for message  $W_{k-1}$ ) in the K - k + 2 messages setting to the interference in the K - k + 1 messages setting. Note that Lemma 1 is a generalization of the corresponding lemma in the single round PIR setting (see [6, Lemma 5]). In the proof, the intuitive idea is to reduce the interference contained in answers from all N databases to that from T databases and then bound the interference using the privacy and correctness constraints. Note that in the multi-round setting, when expanding the interference term, the causality constraint must not be violated.

Lemma 2:

$$I(W_{2:K}; Q_{1:N}^{[1]}(1:\Gamma), A_{1:N}^{[1]}(1:\Gamma), \mathbb{F}|W_1, \mathbb{G}) < L(1/R-1) + o(L).$$
(45)

<sup>&</sup>lt;sup>10</sup>Since optimal compression limits are typically achieved asymptotically, if the data is not assumed to be uniform a-priori, then as noted by [20] and [21] the privacy guarantees would also be subject to  $\epsilon$ -leakage that approaches zero as message length approaches infinity.

|            | Prob. 1/2             |                       | Prob. 1/2             |                       |
|------------|-----------------------|-----------------------|-----------------------|-----------------------|
|            | Want $W_1$            | Want $W_2$            | Want $W_1$            | Want $W_2$            |
| Database 1 | $a_1, b_1, a_2 + b_2$ | $a_1, b_1, a_2 + b_2$ | $a_3, b_3, a_4 + b_4$ | $a_3, b_3, a_4 + b_4$ |
| Database 2 | $a_4, b_2, a_3 + b_1$ | $a_2, b_4, a_1 + b_3$ | $a_2, b_4, a_1 + b_3$ | $a_4, b_2, a_3 + b_1$ |

$$NI(W_{k:K}; \mathcal{Q}_{1:N}^{[k-1]}(1:\Gamma), A_{1:N}^{[k-1]}(1:\Gamma), \mathbb{F}|W_{1:k-1}, \mathbb{G}) \\ \geq \frac{N}{\binom{N}{T}} \sum_{\mathcal{T} \subset [1:N], |\mathcal{T}|=T} I(W_{k:K}; \mathcal{Q}_{\mathcal{T}}^{[k-1]}(1:\Gamma), A_{\mathcal{T}}^{[k-1]}(1:\Gamma)|W_{1:k-1}, \mathbb{G})$$
(32)

$$\stackrel{(11)}{=} \frac{N}{\binom{N}{T}} \sum_{\mathcal{T} \subset [1:N], |\mathcal{T}|=T} I(W_{k:K}; \mathcal{Q}_{\mathcal{T}}^{[k]}(1:\Gamma), A_{\mathcal{T}}^{[k]}(1:\Gamma)|W_{1:k-1}, \mathbb{G})$$
(33)

$$= \frac{N}{\binom{N}{T}} \sum_{\mathcal{T} \subset [1:N], |\mathcal{T}|=T} \sum_{\gamma=1}^{\Gamma} I(W_{k:K}; Q_{\mathcal{T}}^{[k]}(\gamma), A_{\mathcal{T}}^{[k]}(\gamma)|Q_{\mathcal{T}}^{[k]}(1:\gamma-1), A_{\mathcal{T}}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{G})$$

$$\geq \frac{N}{\binom{N}{T}} \sum_{\mathcal{T}} \sum_{T} I(W_{k:K}; A_{\mathcal{T}}^{[k]}(\gamma)|Q_{\mathcal{T}}^{[k]}(1:\gamma), A_{\mathcal{T}}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{G})$$
(34)

$$\binom{N}{T}_{\mathcal{T}\subset [1:N], |\mathcal{T}|=T} \sum_{\gamma=1}^{\Gamma} H(A_{\mathcal{T}}^{[k]}(\gamma)|O_{1,\gamma}^{[k]}(1:\gamma), A_{1,\gamma}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{F}, \mathbb{G})$$
(36)

$$\frac{\overline{\binom{N}{T}}}{\prod_{\mathcal{T}\subset [1:N], |\mathcal{T}|=T}} \sum_{\gamma=1}^{\mathcal{T}} \frac{\mathcal{T}(A_{\mathcal{T}}(\gamma)|\mathcal{Q}_{1:N}(1:\gamma), A_{1:N}(1:\gamma-1), w_{1:k-1}, \mathbb{F}, \mathbb{G})}{\Gamma}$$
(50)

$$\geq T \sum_{\gamma=1}^{\infty} H(A_{1:N}^{[k]}(\gamma) | Q_{1:N}^{[k]}(1:\gamma), A_{1:N}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{F}, \mathbb{G}) \quad (\text{Han's inequality [30]})$$
(37)

$$\stackrel{(8)(10)}{=} T \sum_{\gamma=1}^{1} I(W_{k:K}; A_{1:N}^{[k]}(\gamma) | \mathcal{Q}_{1:N}^{[k]}(1:\gamma), A_{1:N}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{F}, \mathbb{G})$$
(38)

$$\stackrel{(7)(9)}{=} T \sum_{\gamma=1}^{1} I(W_{k:K}; \mathcal{Q}_{1:N}^{[k]}(\gamma), A_{1:N}^{[k]}(\gamma) | \mathcal{Q}_{1:N}^{[k]}(1:\gamma-1), A_{1:N}^{[k]}(1:\gamma-1), W_{1:k-1}, \mathbb{F}, \mathbb{G})$$
(39)

$$= TI(W_{k:K}; Q_{1:N}^{[k]}(1:\Gamma), A_{1:N}^{[k]}(1:\Gamma)|W_{1:k-1}, \mathbb{F}, \mathbb{G})$$
<sup>(12)</sup>
<sup>(13)</sup>
<sup>(13)</sup>
<sup>(13)</sup>
<sup>(14)</sup>

$$\stackrel{(13)}{=} TI(W_{k:K}; W_k, Q_{1:N}^{[k]}(1:\Gamma), A_{1:N}^{[k]}(1:\Gamma) | W_{1:k-1}, \mathbb{F}, \mathbb{G}) - o(L)T$$
(41)

$$TI(W_{k:K}; W_k | W_{1:k-1}, \mathbb{F}, \mathbb{G}) - o(L)T$$

$$+ TI(W_{k+1:K}; Q_{1:N}^{K_{1}}(1:\Gamma), A_{1:N}^{K_{1}}(1:\Gamma)|W_{1:k}, \mathbb{F}, \mathbb{G})$$

$$(42)$$

$$(3)(6) LT = To(L) + TL(W_{k-1} = O^{[k]}(1:\Gamma) + A^{[k]}(1:\Gamma)|W_{k-1} = \mathbb{F} |\Omega|$$

$$(42)$$

$$\stackrel{5}{=} LT - To(L) + TI(W_{k+1:K}; Q_{1:N}^{[k]}(1:\Gamma), A_{1:N}^{[k]}(1:\Gamma)|W_{1:k}, \mathbb{F}, \mathbb{G})$$
(43)

$$\stackrel{(6)}{=} LT - To(L) + TI(W_{k+1:K}; Q_{1:N}^{[k]}(1:\Gamma), A_{1:N}^{[k]}(1:\Gamma), \mathbb{F}|W_{1:k}, \mathbb{G})$$
(44)

Proof:

=

$$I(W_{2:K}; Q_{1:N}^{[1]}(1:\Gamma), A_{1:N}^{[1]}(1:\Gamma), \mathbb{F}|W_1, \mathbb{G}) \stackrel{(6)}{=} I(W_{2:K}; Q_{1:N}^{[1]}(1:\Gamma), A_{1:N}^{[1]}(1:\Gamma), W_1, \mathbb{F}, \mathbb{G})$$
(46)  
$$\stackrel{(7)(9)}{=} I(W_{2:K}; A_{1:N}^{[1]}(1:\Gamma), W_1, \mathbb{F}, \mathbb{G})$$
(47)

$$= I(W_{2:K}; A_{1:N}^{[1]}(1:\Gamma), \mathbb{F}, \mathbb{G}) + I(W_{2:K}; W_1 | A_{1:N}^{[1]}(1:\Gamma), \mathbb{F}, \mathbb{G})$$
(48)

$$\stackrel{(6)(13)}{=} I(W_{2:K}; A^{[1]}_{1:N}(1:\Gamma) | \mathbb{F}, \mathbb{G}) + o(L)$$

$$= H(A^{[1]}_{1:N}(1:\Gamma) | \mathbb{F}, \mathbb{G})$$

$$(49)$$

$$- H(A_{1:N}^{[1]}(1:\Gamma)|\mathbb{F}, \mathbb{G}, W_{2:K}) + o(L)$$
(50)

$$\stackrel{(12)}{\leq} L/R - H(A_{1:N}^{[1]}(1:\Gamma)|\mathbb{F}, \mathbb{G}, W_{2:K}) + o(L)$$
(51)  
$$\stackrel{(13)}{=} L/R - H(W_1, A_{1:N}^{[1]}(1:\Gamma)|\mathbb{F}, \mathbb{G}, W_{2:K}) + o(L)$$

$$\leq L/R - H(W_1|\mathbb{F}, \mathbb{G}, W_{2:K}) + o(L)$$
(53)

$$\stackrel{(6)}{=} L/R - L + o(L)L = L(1/R - 1) + o(L)$$
(54)

*Remark:* The intuition of Lemma 2 is that among the total download (L/R symbols), to leave L symbols for the desired message, the interference about all other messages must have size no more than L/R - L symbols.

With Lemma 1 and Lemma 2, we are ready to prove the converse.

*Rate Outerbound:* Starting from k = 2 and applying (31) repeatedly for  $k \in [3: K]$ ,

$$I(W_{2:K}; Q_{1:N}^{[1]}(1:\Gamma), A_{1:N}^{[1]}(1:\Gamma), \mathbb{F}|W_{1}, \mathbb{G})$$

$$\stackrel{(31)}{\geq} \frac{T}{N}I(W_{3:K}; Q_{1:N}^{[2]}(1:\Gamma), A_{1:N}^{[2]}(1:\Gamma), \mathbb{F}|W_{1}, W_{2}, \mathbb{G})$$

$$+ \frac{LT(1 - o(L)/L)}{N}$$

$$\stackrel{(31)}{\geq} \cdots \qquad (55)$$

$$\stackrel{(31)}{\geq} \frac{T^{K-2}}{N^{K-2}}I(W_{K}; Q_{1:N}^{[K-1]}(1:\Gamma), A_{1:N}^{[K-1]}(1:\Gamma), \mathbb{F}| \dots$$

$$\dots W_{1:K-1}, \mathbb{G})$$

$$+ \frac{LT(1 - o(L)/L)}{N} + \dots + \frac{LT^{K-2}(1 - o(L)/L)}{N^{K-2}}$$

$$\stackrel{(31)}{\geq} \frac{T^{K-2}}{N^{K-2}}\frac{LT(1 - o(L)/L)}{N} + \frac{LT(1 - o(L)/L)}{N}$$

$$+ \dots + \frac{LT^{K-2}(1 - o(L)/L)}{N^{K-2}} \qquad (56)$$

$$(L - (L))(T)(M + U - V_{K-1})(M^{K-1})$$

$$= (L - o(L))(T/N + \dots + T^{K-1}/N^{K-1})$$
(57)

Combining (57) and (45), we have

$$L(1/R - 1) + o(L) \ge (L - o(L))(T/N + \dots + T^{K-1}/N^{K-1})$$
(58)

Normalizing by L and letting L go to infinity gives us

$$1/R - 1 \ge T/N + \dots + T^{K-1}/N^{K-1}$$
(59)

$$\Rightarrow R \le (1 + T/N + \dots + T^{K-1}/N^{K-1})^{-1}$$
 (60)

thus, the proof is complete.

## V. PROOF OF THEOREM 2 – STATEMENT 2.

We show that when K = 2, N = 2, T = 1,  $\Gamma = 1$  and the rate equals 2/3, the storage overhead of all zero-error,<sup>11</sup> linear, and single-round PIR schemes is no less than 3/2. Since we only consider single-round schemes in this section, we will simplify the notation, e.g., instead of  $Q_2^{[1]}(1)$  we write simply  $Q_2^{[1]}$ . In addition, without loss of generality, let us make the following simplifying assumptions.

1) We assume that the PIR scheme is symmetric, in that

$$H(A_1^{[1]}|\mathbb{F},\mathbb{G}) = H(A_2^{[1]}|\mathbb{F},\mathbb{G}) = H(A_2^{[2]}|\mathbb{F},\mathbb{G}) \quad (61)$$
$$H(S_1) = H(S_2) \quad (62)$$

Given any (asymmetric) PIR scheme that retrieves messages of size L, a symmetric PIR scheme with the same rate and storage overhead that retrieves messages of size NL is obtained by repeating the original scheme N times, and in the  $n^{th}$  repetition shifting the database indices cyclically by n. This symmetrization process is described in Lemma 4 (see Section V-A).

2) We assume that  $Q_1^{[1]} = Q_1^{[2]}$ , i.e., the query for the first database is chosen without the knowledge of the desired message index. There is no loss of generality in

this assumption because of the privacy constraint, which requires that  $Q_1^{[\theta]}$  is independent of  $\theta$ .<sup>12</sup> Note that this also means that  $A_1^{[1]} = A_1^{[2]}$ .

Our goal is to prove a lower bound on the storage overhead. Since the PIR scheme is symmetric by assumption, the storage overhead is  $(H(S_1) + H(S_2))/2L = H(S_2)/L$ . Furthermore,  $H(S_2) \ge H(A_2^{[1]}, A_2^{[2]}|\mathbb{F}, \mathbb{G})$ , so we will prove a lower bound on  $H(A_2^{[1]}, A_2^{[2]}|\mathbb{F}, \mathbb{G})$  instead.

Let us start with a useful lemma that holds for all linear and non-linear schemes.

Lemma 3:

$$H(A_{1}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G}) = H(A_{2}^{[2]}|W_{1}, \mathbb{F}, \mathbb{G})$$
  

$$= H(A_{2}^{[2]}|W_{2}, \mathbb{F}, \mathbb{G})$$
  

$$= L/2 \qquad (63)$$
  

$$H(A_{2}^{[2]}|W_{1}, A_{2}^{[1]}, \mathbb{F}, \mathbb{G}) = H(A_{2}^{[2]}|W_{2}, A_{2}^{[1]}, \mathbb{F}, \mathbb{G})$$
  

$$= L/2 \qquad (64)$$

*Proof:* We prove (63) first. On the one hand, we substitute<sup>13</sup> R = 2/3 in Lemma 2. Then from (47) - (54), we have

$$L/2 \ge I(W_2; A_1^{[1]}, A_2^{[1]}, W_1, \mathbb{F}, \mathbb{G})$$
 (65)

$$\stackrel{(6)}{=} I(W_2; A_1^{[1]}, A_2^{[1]} | W_1, \mathbb{F}, \mathbb{G})$$
(66)

$$\stackrel{(7)(8)(4)}{=} H(A_1^{[1]}, A_2^{[1]} | W_1, \mathbb{F}, \mathbb{G})$$
(67)

$$\Rightarrow L/2 \ge H(A_1^{[1]}|W_1, \mathbb{F}, \mathbb{G})$$
(68)

and

$$L/2 \ge H(A_2^{[1]}|W_1, \mathbb{F}, \mathbb{G}) \tag{69}$$

On the other hand, from (32) - (44), as shown at the top of the previous page, in Lemma 1, we have

$$L \leq I(W_2; Q_1^{[1]}, A_1^{[1]} | W_1, \mathbb{G}) + I(W_2; Q_2^{[1]}, A_2^{[1]} | W_1, \mathbb{G})$$
(70)

$$\leq I(W_2; Q_1^{[1]}, A_1^{[1]}, \mathbb{F}|W_1, \mathbb{G}) + I(W_2; Q_2^{[1]}, A_2^{[1]}, \mathbb{F}|W_1, \mathbb{G})$$
(71)

$$\stackrel{(6)}{=} I(W_2; Q_1^{[1]}, A_1^{[1]} | W_1, \mathbb{F}, \mathbb{G}) + I(W_2; Q_2^{[1]}, A_2^{[1]} | W_1, \mathbb{F}, \mathbb{G})$$
(72)

$$\stackrel{(7)(8)(4)}{=} H(A_1^{[1]}|W_1, \mathbb{F}, \mathbb{G}) + H(A_2^{[1]}|W_1, \mathbb{F}, \mathbb{G}) \quad (73)$$

Combining (68), (69) and (73), we have shown that

$$H(A_1^{[1]}|W_1, \mathbb{F}, \mathbb{G}) = H(A_2^{[1]}|W_1, \mathbb{F}, \mathbb{G}) = L/2$$
(74)

Symmetrically, it follows that  $H(A_2^{[2]}|W_2, \mathbb{F}, \mathbb{G}) = L/2$ . We are left to prove  $H(A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) = L/2$ . On the one hand,

<sup>12</sup>Note that instead of  $Q_1^{[1]} = Q_1^{[2]}$ , we could equivalently assume that  $Q_2^{[1]} = Q_2^{[2]}$  without of loss of generality (because privacy also requires that  $Q_2^{[\theta]}$  is independent of  $\theta$ ). However, if we simultaneously assume both  $Q_1^{[1]} = Q_1^{[2]}$  and  $Q_2^{[1]} = Q_2^{[2]}$ , then there is a loss of generality because together  $(Q_1^{[\theta]}, Q_2^{[\theta]})$  is *not* required to be independent of  $\theta$  by the privacy constraint.

<sup>&</sup>lt;sup>11</sup>Our converse proof extends to the  $\epsilon$ -error case.

<sup>&</sup>lt;sup>13</sup>Since we are considering only zero-error schemes, the o(L) term in Lemma 2 is exactly 0.

from (68) and (69), we have

$$L/2 \geq H(A_1^{[1]}|W_1, \mathbb{F}, \mathbb{G}) = H(A_1^{[2]}|W_1, \mathbb{F}, \mathbb{G}) \quad (\text{Using } A_1^{[1]} = A_1^{[2]}) \quad (75)$$

$$L/2 \geq H(A_2^{[1]}|W_1, \mathbb{F}, \mathbb{G})$$
(76)

$$\stackrel{(')}{=} H(A_2^{[1]}|W_1, Q_2^{[1]}, \mathbb{F}, \mathbb{G})$$
(77)

$$\stackrel{(6)}{=} H(A_2^{[1]}|W_1, Q_2^{[1]}, \mathbb{G})$$
(78)

$$\stackrel{(11)}{=} H(A_2^{[2]}|W_1, Q_2^{[2]}, \mathbb{G}) \tag{79}$$

$$\stackrel{(8)(6)}{=} H(A_2^{[2]}|W_1, Q_2^{[2]}, \mathbb{F}, \mathbb{G})$$
(80)

$$\stackrel{(I)}{=} H(A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) \tag{81}$$

where (79) follows from the fact that for single-round PIR, the desired message index is independent of the messages, queries and answers. Its detailed proof is presented as follows. Note that<sup>14</sup>

$$I(Q_2^{[\theta]}, \theta, \mathbb{F}; W_1, W_2, \mathbb{G}) \stackrel{(7)}{=} I(\theta, \mathbb{F}; W_1, W_2, \mathbb{G})$$

$$\stackrel{(6)}{=} 0 \qquad (82)$$

$$\implies I(Q_2^{[\theta]}; W_1, W_2, \mathbb{G}) = I(Q_2^{[\theta]}; W_1, W_2, \mathbb{G}|\theta)$$

$$= 0 \qquad (83)$$

Next,

$$I(\theta; W_1, W_2, \mathbb{G}, Q_2^{[\theta]}) = I(\theta; W_1, W_2, \mathbb{G}) + I(\theta; Q_2^{[\theta]} | W_1, W_2, \mathbb{G})$$
(84)

$$\stackrel{(0)}{=} I(\theta; Q_2^{[\theta]} | W_1, W_2, \mathbb{G})$$
(85)

$$= H(Q_2^{[\theta]}|W_1, W_2, \mathbb{G}) - H(Q_2^{[\theta]}|\theta, W_1, W_2, \mathbb{G})$$
(86)

$$\stackrel{(83)}{=} H(Q_2^{[\theta]}) - H(Q_2^{[\theta]}|\theta) \tag{87}$$

$$\stackrel{(11)}{=} 0 \tag{88}$$

$$\Longrightarrow W_1, W_2, Q_2^{[1]}, \mathbb{G} \sim W_1, W_2, Q_2^{[2]}, \mathbb{G}$$
(89)

$$\stackrel{8)(4)}{\Longrightarrow} A_2^{[1]}, W_1, W_2, Q_2^{[1]}, \mathbb{G} \sim A_2^{[2]}, W_1, W_2, Q_2^{[2]}, \mathbb{G}$$
(90)

$$\implies A_2^{[1]}, W_1, Q_2^{[1]}, \mathbb{G} \sim A_2^{[2]}, W_1, Q_2^{[2]}, \mathbb{G}$$
(91)

(78) and (80) are due to the Markov chain  $\mathbb{F} - (W_1, Q_2^{[k]}, \mathbb{G}) - A_2^{[k]}, k = 1, 2$ , which is proved as follows.

$$I(A_{2}^{[k]}; \mathbb{F}|W_{1}, Q_{2}^{[k]}, \mathbb{G}) \\\leq I(A_{2}^{[k]}, S_{2}; \mathbb{F}|W_{1}, Q_{2}^{[k]}, \mathbb{G})$$
(92)

$$= I(S_2; \mathbb{F}|W_1, Q_2^{[k]}, \mathbb{G}) + I(A_2^{[k]}; \mathbb{F}|W_1, Q_2^{[k]}, \mathbb{G}, S_2)$$
(93)

$$\stackrel{(8)}{=} I(S_2; \mathbb{F}|W_1, Q_2^{[k]}, \mathbb{G}) \tag{94}$$

$$\leq I(S_2, W_2; \mathbb{F}|W_1, Q_2^{[k]}, \mathbb{G})$$
(95)

$$= I(W_2; \mathbb{F}|W_1, Q_2^{[K]}, \mathbb{G})$$

<sup>14</sup>The distribution of  $Q_2^{[\theta]}$  is a mixture of the distributions of  $Q_2^{[1]}$  and  $Q_2^{[2]}$ . Conditioned on  $\theta = 1$ ,  $Q_2^{[\theta]} = Q_2^{[1]}$ . Conditioned on  $\theta = 2$ .  $Q_2^{[\theta]} = Q_2^{[2]}$ . The privacy condition (11) can be equivalently expressed as  $I(\theta; Q_n^{[\theta]}, A_n^{[\theta]}, \mathbb{G}, S_n) = 0, n \in \{1, 2\}$ , in this case.

$$+ I(S_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G}, W_1, W_2)$$
 (96)

$$\stackrel{(4)}{\leq} I(W_2; \mathbb{F}, W_1, Q_2^{[k]}, \mathbb{G}) \tag{97}$$

$$\stackrel{(7)(6)}{=} 0$$
 (98)

On the other hand, from (70), we have

$$L \leq I(W_2; Q_1^{[1]}, A_1^{[1]}|W_1, \mathbb{G}) + I(W_2; Q_2^{[1]}, A_2^{[1]}|W_1, \mathbb{G})$$
<sup>(11)</sup>

$$\stackrel{\cong}{=} I(W_2; Q_1^{[2]}, A_1^{[2]}|W_1, \mathbb{G}) + I(W_2; Q_2^{[2]}, A_2^{[2]}|W_1, \mathbb{G})$$
(100)

$$\leq I(W_2; Q_1^{[2]}, A_1^{[2]}, \mathbb{F}|W_1, \mathbb{G}) + I(W_2; Q_1^{[2]}, A_1^{[2]}, \mathbb{F}|W_1, \mathbb{G})$$
(101)

$$\stackrel{(6)}{=} I(W_2; Q_1^{[2]}, A_1^{[2]} | W_1, \mathbb{F}, \mathbb{G})$$
(101)

$$+ I(W_2; Q_2^{[2]}, A_2^{[2]} | W_1, \mathbb{F}, \mathbb{G})$$
(102)  
(7)(8)(4)  $W(A^{[2]} | W_1, \mathbb{F}, \mathbb{G}) = W(A^{[2]} | W_1, \mathbb{F}, \mathbb{G})$ (102)

$$\stackrel{(1)}{=} H(A_1^{[2]}|W_1, \mathbb{F}, \mathbb{G}) + H(A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) \quad (103)$$

Combining (75), (81) and (103), we have shown that  $H(A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) = L/2$ . The proof of (63) is complete. Next we prove (64). On the one hand,

$$H(A_2^{[2]}|W_1, A_2^{[1]}, \mathbb{F}, \mathbb{G}) \le H(A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) \stackrel{(63)}{=} L/2 \quad (104)$$

On the other hand, from sub-modularity of entropy functions we have

$$H(A_{2}^{[2]}, A_{2}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G})$$

$$\geq -H(A_{2}^{[1]}, A_{1}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G})$$

$$+ H(A_{1}^{[1]}, A_{2}^{[2]}, A_{2}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G})$$

$$+ H(A_{2}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G})$$

$$\stackrel{(67)(13)(74)}{\geq} -L/2 + H(A_{1}^{[1]}, A_{2}^{[2]}, A_{2}^{[1]}, W_{2}|W_{1}, \mathbb{F}, \mathbb{G})$$

$$(105)$$

$$+L/2$$
 (106)

$$\geq H(W_2|W_1, \mathbb{F}, \mathbb{G}) \stackrel{(6)}{=} L \tag{107}$$

$$\Rightarrow H(A_{2}^{[2]}|W_{1}, A_{2}^{[1]}, \mathbb{F}, \mathbb{G}) = H(A_{2}^{[2]}, A_{2}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G}) - H(A_{2}^{[1]}|W_{1}, \mathbb{F}, \mathbb{G}) \stackrel{(74)}{\geq} L/2$$
(108)

Note that the second term of (106) follows from the assumption that  $A_1^{[1]} = A_1^{[2]}$  so that from  $A_1^{[1]}, A_2^{[2]}$ , we can decode  $W_2$  just as from  $A_1^{[2]}, A_2^{[2]}$ , we can decode  $W_2$ . Combining (104), (108), we have proved  $H(A_2^{[2]}|W_1, A_2^{[1]}, \mathbb{F}, \mathbb{G}) = L/2$ . Symmetrically, it follows that  $H(A_2^{[2]}|W_2, A_2^{[1]}, \mathbb{F}, \mathbb{G}) = L/2$ . Therefore, the desired inequality (64) is obtained.

To proceed, we need Ingleton's inequality, which is stated as follows.

Theorem 3 (Ingleton's Inequality [31]): For four subspaces of a given finite vector space, A, B, C, D, the following

inequality holds.15

$$I(A; B) \le I(A; B|C) + I(A; B|D) + I(C; D) \quad (109)$$

For a given value of  $\mathbb{F}$ ,  $\mathbb{G}$ ,  $A_2^{[1]}$ ,  $A_2^{[2]}$  are linear combinations of  $W_1$ ,  $W_2$  with constant coefficients as we consider linear schemes. So we set  $A = A_2^{[1]}$ ,  $B = A_2^{[2]}$ ,  $C = W_1$ ,  $D = W_2$  (for given  $\mathbb{F}$ ,  $\mathbb{G}$ ). Note that from Lemma 3, we know that  $I(A_2^{[1]}; A_2^{[2]}|W_1, \mathbb{F}, \mathbb{G}) = I(A_2^{[1]}; A_2^{[2]}|W_2, \mathbb{F}, \mathbb{G}) = 0$ . Plugging in (109) that holds for linear schemes but not for non-linear schemes, we have

$$I(A_{2}^{[1]}; A_{2}^{[2]} | \mathbb{F}, \mathbb{G}) \\ \leq I(A_{2}^{[1]}; A_{2}^{[2]} | W_{1}, \mathbb{F}, \mathbb{G}) + I(A_{2}^{[1]}; A_{2}^{[2]} | W_{2}, \mathbb{F}, \mathbb{G}) \\ + \underbrace{I(W_{1}; W_{2} | \mathbb{F}, \mathbb{G})}_{=0 \text{ from } (0)} = 0$$
(110)

$$\Longrightarrow H(A_2^{[1]}, A_2^{[2]} | \mathbb{F}, \mathbb{G})$$

$$= H(A_2^{[1]} | \mathbb{F}, \mathbb{G}) + H(A_2^{[2]} | \mathbb{F}, \mathbb{G})$$

$$(111)$$

$$\stackrel{(61)}{=} H(A_2^{[1]}|\mathbb{F}, \mathbb{G}) + H(A_1^{[1]}|\mathbb{F}, \mathbb{G})$$
(112)

$$\geq H(A_1^{[1]}, A_2^{[1]} | \mathbb{F}, \mathbb{G})$$
(113)

$$\stackrel{(13)}{=} H(W_1, A_1^{[1]}, A_2^{[1]} | \mathbb{F}, \mathbb{G})$$
(114)

$$= H(W_1|\mathbb{F}, \mathbb{G}) + H(A_1^{[1]}, A_2^{[1]}|W_1, \mathbb{F}, \mathbb{G})$$
(115)

$$\stackrel{(6)}{\geq} L + H(A_1^{[1]}|W_1, \mathbb{F}, \mathbb{G}) \stackrel{(63)}{=} 3L/2 \tag{116}$$
$$\implies \alpha = H(S_2)/L$$

$$\geq H(A_2^{[1]}, A_2^{[2]} | \mathbb{F}, \mathbb{G}) / L \geq 3/2$$
(117)

# A. Symmetrization

*Lemma 4:* <sup>16</sup>Consider the single-round PIR problem with K = 2 messages and N = 2databases. Suppose we have a scheme described by  $\overline{L}, \overline{W}_1, \overline{W}_2, \overline{S}_1, \overline{S}_2, \overline{Q}_{1:2}^{[1]}, \overline{Q}_{1:2}^{[2]}, \overline{A}_{1:2}^{[1]}, \overline{F}, \overline{\mathbb{G}}$ . Then we can construct a symmetric PIR scheme, also for K = N = 2, described by  $L, W_1, W_2, S_1, S_2, Q_{1:2}^{[1]}, Q_{1:2}^{[2]}, A_{1:2}^{[1]}, A_{1:2}^{[2]}, \mathbb{F}, \mathbb{G}$ such that

$$H(A_1^{[1]}|\mathbb{F},\mathbb{G}) = H(A_2^{[1]}|\mathbb{F},\mathbb{G}) = H(A_2^{[2]}|\mathbb{F},\mathbb{G}) \quad (118)$$

$$H(S_1) = H(S_2)$$
 (119)

$$L = 2\bar{L} \tag{120}$$

such that the symmetric PIR scheme has the same rate and storage overhead as the original PIR scheme.

*Proof:* Consider two independent implementations of the asymmetric PIR scheme. Let us use the 'bar' notation for the first implementation and the 'tilde' notation for the second implementation. In the first implementation, there are two

messages  $\overline{W}_1$ ,  $\overline{W}_2$ , each of length  $\overline{L}$ , two databases DB1 and DB2 which store  $\overline{S}_1$ ,  $\overline{S}_2$ , respectively. In the second implementation, there are two messages  $\widetilde{W}_1$ ,  $\widetilde{W}_2$ , each of length  $\widetilde{L} = \overline{L}$ , two databases DB2 and DB1 which store  $\widetilde{S}_1$ ,  $\widetilde{S}_2$ , respectively. Note the critical detail that the database indices are switched in the second implementation. The asymmetric PIR scheme specifies the queries for each implementation such that the user can privately retrieve an arbitrarily chosen message from each implementation.

The symmetric PIR scheme is created by combining the two implementations. In the combined scheme, there are two messages  $W_1 = (\bar{W}_1, \tilde{W}_1)$  and  $W_2 = (\bar{W}_2, \tilde{W}_2)$ , each of length  $L = 2\bar{L}$ , two databases DB1 and DB2 which store  $(\bar{S}_1, \bar{S}_2)$  and  $(\bar{S}_2, \bar{S}_1)$ , respectively. Retrieval works exactly as before. For example, if the user wishes to privately retrieve  $W_1 = (\bar{W}_1, \tilde{W}_1)$ , then it retrieves  $\bar{W}_1$  exactly as in the first implementation, and  $\tilde{W}_1$  exactly as in the second implementation.

Since the symmetric scheme is comprised of two independent implementations of the original PIR scheme, the message size, total download size, total storage size, are all doubled relative to the original PIR scheme. As a result the rate and storage overhead, which are normalized quantities, remain unchanged in the new scheme. Symmetry is achieved because each database from the original PIR scheme is equally represented within each database in the new PIR scheme.

Mathematically,

$$W_1 = (\bar{W}_1, \tilde{W}_1), W_2 = (\bar{W}_2, \tilde{W}_2)$$
 (121)

$$S_1 = (\bar{S}_1, \tilde{S}_2), S_2 = (\bar{S}_2, \tilde{S}_1)$$
 (122)

$$\mathbb{F} = (\mathbb{F}, \mathbb{F}), \mathbb{G} = (\mathbb{G}, \mathbb{G})$$
(123)

$$Q_n^{[k]} = (\tilde{Q}_n^{[k]}, \tilde{Q}_{3-n}^{[k]}), \quad n = 1, 2, \ k = 1, 2$$
(124)

$$A_n^{[k]} = (A_n^{[k]}, \tilde{A}_{3-n}^{[k]})$$
(125)

where each random variable with a bar symbol is independent of and identically distributed with the same random variable with a tilde symbol. We are now ready to prove the first equality in (118).

$$H(A_1^{[1]}|\mathbb{F},\mathbb{G}) = H(\bar{A}_1^{[1]},\tilde{A}_2^{[1]}|\mathbb{F},\mathbb{G})$$
(126)

$$= H(\bar{A}_{1}^{[1]}|\bar{\mathbb{F}},\bar{\mathbb{G}}) + H(\tilde{A}_{2}^{[1]}|\tilde{\mathbb{F}},\bar{\mathbb{G}}) \qquad (127)$$

$$= H(\tilde{A}_{2}^{[1]}|\tilde{\mathbb{F}}|\tilde{\mathbb{G}}) + H(\tilde{A}_{2}^{[1]}|\tilde{\mathbb{F}}|\tilde{\mathbb{G}}) \quad (128)$$

$$= H(\bar{A}_{1}^{[1]}, \tilde{A}_{1}^{[1]} | \mathbb{F}, \mathbb{G})$$
(129)

$$= H(A_{2}^{[1]}|\mathbb{F},\mathbb{G})$$
(130)

$$H(A_2^{[1]}|\mathbb{F},\mathbb{G}) \tag{130}$$

where (127) and (129) follow from the fact that the two copies of the given scheme are independent and (128) is due to the property that the two copies are identically distributed. Consider the second equality in (118).

$$H(A_2^{[1]}|\mathbb{F},\mathbb{G}) \stackrel{(7)}{=} H(A_2^{[1]}|Q_2^{[1]},\mathbb{F},\mathbb{G})$$
(131)

$$= H(A_2^{[1]}|Q_2^{[1]}, \mathbb{G})$$
(132)

$$\stackrel{(11)}{=} H(A_2^{[2]}|Q_2^{[2]}, \mathbb{G}) \tag{133}$$

$$= H(A_{2}^{[2]}|Q_{2}^{[2]}, \mathbb{F}, \mathbb{G})$$
(134)

$$\stackrel{(f)}{=} H(A_2^{[2]}|\mathbb{F},\mathbb{G}) \tag{135}$$

<sup>&</sup>lt;sup>15</sup>For subspaces A, B, we follow the convention that H(A) represents the dimension of subspace A, and H(A, B) represents the dimension of the vector space spanned by the union of the subspaces A, B. Using this convention, inequalities on the dimensions of subspaces can be expressed using information theoretic measures, such as (joint) entropy and (conditional) mutual information. Ingleton's inequality has been stated in this form in prior work in information theory literature, e.g., [32], [33].

<sup>&</sup>lt;sup>16</sup>Extensions of this symmetrization lemma to multiple rounds, arbitrary number of messages and databases may be similarly obtained.

where (132) and (134) are due to the Markov chain  $\mathbb{F} - (Q_2^{[k]}, \mathbb{G}) - A_2^{[k]}, k = 1, 2$ , which is proved as follows.

$$I(A_{2}^{[k]}; \mathbb{F}|Q_{2}^{[k]}, \mathbb{G}) \\\leq I(A_{2}^{[k]}, S_{2}; \mathbb{F}|Q_{2}^{[k]}, \mathbb{G})$$
(136)

$$= I(S_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G}) + I(A_2^{[k]}; \mathbb{F}|Q_2^{[k]}, \mathbb{G}, S_2) \quad (137)$$

$$\stackrel{(8)}{=} I(S_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G}) \tag{138}$$

$$\leq I(S_2, W_1, W_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G})$$
(139)

$$= I(W_1, W_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G})$$

$$+ I(S_2; \mathbb{F}|Q_2^{[k]}, \mathbb{G}, W_1, W_2)$$
 (140)

$$\stackrel{(4)}{\leq} I(W_1, W_2; \mathbb{F}, Q_2^{[k]}, \mathbb{G})$$
(141)

$$\stackrel{(7)(6)}{=} 0 \tag{142}$$

Finally, we prove (119).

(1)

$$H(S_1) = H(\bar{S}_1, \tilde{S}_2)$$
 (143)

$$= H(\bar{S}_1) + H(\tilde{S}_2)$$
(144)

$$= H(\tilde{S}_1) + H(\bar{S}_2) \tag{145}$$

$$= H(\bar{S}_2, \tilde{S}_1) \tag{146}$$

$$=H(S_2) \tag{147}$$

where (144) and (146) follow from the fact that the two copies of the given scheme are independent and (145) is due to the property that the two copies are identically distributed.

## VI. CONCLUSION

We showed that the capacity of MPIR is equal to the capacity of PIR, both with and without T-privacy constraints. Our result implies that there is no advantage in terms of capacity from multiround over single-round schemes, nonlinear over linear schemes, or  $\epsilon$ -error over zero-error schemes. We also offered a counterpoint to this pessimistic result by exploring optimal storage overhead instead of capacity. Specifically, we constructed a simple multiround, non-linear,  $\epsilon$ -error PIR scheme that achieves a strictly smaller storage overhead than the best possible with any single-round, linear, zero-error PIR scheme. The simplicity of the scheme makes it an attractive point of reference for future work toward understanding the role of linear versus non-linear schemes, zeroerror versus  $\epsilon$ -error capacity, and single-round versus multiple round communications. Another interesting insight revealed by the scheme is the privacy benefit of reduced storage overhead. By not storing all the information at each database, and by optimally compressing the stored information, not only do we reduce the storage overhead, but also we enable privacy where it wouldn't hold otherwise.

#### REFERENCES

- B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan, "Private information retrieval," in *Proc. 36th Annu. Symp. Found. Comput. Sci.*, 1995, pp. 41–50.
- [2] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, "Private information retrieval," J. ACM, vol. 45, no. 6, pp. 965–981, 1998.
- [3] R. Ostrovsky and W. E. Skeith, III, "A survey of single-database private information retrieval: Techniques and applications," in *Public Key Cryptography—PKC*. Berlin, Germany: Springer, 2007, pp. 393–411.

- [4] W. Gasarch, "A survey on private information retrieval," *Bull. EATCS*, vol. 82, pp. 72–107, Feb. 2004.
- [5] S. Yekhanin, "Private information retrieval," *Commun. ACM*, vol. 53, no. 4, pp. 68–73, 2010.
- [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] H. Sun and S. A. Jafar, "Optimal download cost of private information retrieval for arbitrary message length," *IEEE Trans. Inf. Forensics Security*, vol. 12, no. 12, pp. 2920–2932, Dec. 2017.
- [8] H. Sun and S. A. Jafar. (2016). "The capacity of robust private information retrieval with colluding databases." [Online]. Available: https://arxiv.org/abs/1605.00635
- [9] H. Sun and S. A. Jafar. (2016). "The capacity of symmetric private information retrieval." [Online]. Available: https://arxiv.org/ abs/1606.08828
- [10] K. Banawan and S. Ulukus. (2016). "The capacity of private information retrieval from coded databases." [Online]. Available: https://arxiv.org/abs/1609.08138
- [11] Q. Wang and M. Skoglund. (2016). "Symmetric private information retrieval for MDS coded distributed storage." [Online]. Available: https://arxiv.org/abs/1610.04530
- [12] T. H. Chan, S.-W. Ho, and H. Yamamoto, "Private information retrieval for coded storage," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2015, pp. 2842–2846.
- [13] R. Dougherty, C. Freiling, and K. Zeger, "Insufficiency of linear coding in network information flow," *IEEE Trans. Inf. Theory*, vol. 51, no. 8, pp. 2745–2759, Aug. 2005.
- [14] T. Chan and A. Grant, "Dualities between entropy functions and network codes," *IEEE Trans. Inf. Theory*, vol. 54, no. 10, pp. 4470–4487, Oct. 2008.
- [15] S. El Rouayheb, A. Sprintson, and C. Georghiades, "On the index coding problem and its relation to network coding and matroid theory," *IEEE Trans. Inf. Theory*, vol. 56, no. 7, pp. 3187–3195, Jul. 2010.
- [16] A. Blasiak, R. Kleinberg, and E. Lubetzky. (Aug. 2011). "Lexicographic products and the power of non-linear network coding." [Online]. Available: http://arxiv.org/abs/1108.2489
- [17] M. Langberg and M. Effros, "Network coding: Is zero error always possible?" in *Proc. 49th Allerton Conf. Commun., Control Comput.*, 2011, pp. 1478–1485.
- [18] S. Jalali, M. Effros, and T. Ho, "On the impact of a single edge on the network coding capacity," in *Proc. Inf. Theory Appl. Workshop (ITA)*, 2011, pp. 1–5.
- [19] O. Kosut and J. Kliewer, "On the relationship between edge removal and strong converses," in *Proc. Int. Symp. Inf. Theory (ISIT)*, 2016, pp. 1779–1783.
- [20] A. Beimel and Y. Ishai, "Information-theoretic private information retrieval: A unified construction," in *Automata, Languages and Programming.* Berlin, Germany: Springer, 2001, pp. 912–926.
- [21] A. Beimel, Y. Ishai, and E. Kushilevitz, "General constructions for information-theoretic private information retrieval," J. Comput. Syst. Sci., vol. 71, no. 2, pp. 213–247, 2005.
- [22] N. N. Shah, K. V. Rashmi, and K. Ramchandran, "One extra bit of download ensures perfectly private information retrieval," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun./Jul. 2014, pp. 856–860.
- [23] A. Fazeli, A. Vardy, and E. Yaakobi, "Codes for distributed PIR with low storage overhead," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2015, pp. 2852–2856.
- [24] R. Tajeddine, O. W. Gnilke, and S. El Rouayheb. (2016). "Private information retrieval from MDS coded data in distributed storage systems." [Online]. Available: https://arxiv.org/abs/1602.01458
- [25] S. Rao and A. Vardy. (2016). "Lower bound on the redundancy of PIR codes." [Online]. Available: https://arxiv.org/abs/1605.01869
- [26] S. Blackburn and T. Etzion. (2016). "PIR array codes with optimal PIR rate." [Online]. Available: https://arxiv.org/abs/1607.00235
- [27] S. R. Blackburn, T. Etzion, and M. B. Paterson. (2016). "PIR schemes with small download complexity and low storage requirements." [Online]. Available: https://arxiv.org/abs/1609.07027
- [28] Y. Zhang, X. Wang, H. Wei, and G. Ge. (2016). "On private information retrieval array codes." [Online]. Available: https://arxiv.org/ abs/1609.09167
- [29] D. Slepian and J. K. Wolf, "Noiseless coding of correlated information sources," *IEEE Trans. Inf. Theory*, vol. 19, no. 4, pp. 471–480, Jul. 1973.
- [30] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Hoboken, NJ, USA: Wiley, 2006.

5754

- [31] A. W. Ingleton, "Representation of matroids in combinatorial mathematics and its applications," *Combinat. Math. Appl.*, vol. 44, pp. 149–167, Jul. 1971.
- [32] R. Dougherty, C. Freiling, and K. Zeger, "Network coding and matroid theory," *Proc. IEEE*, vol. 99, no. 3, pp. 388–405, Mar. 2011.
- [33] D. Hammer, A. Romashchenko, A. Shen, and N. Vereshchagin, "Inequalities for Shannon entropy and Kolmogorov complexity," *J. Comput. Syst. Sci.*, vol. 60, no. 2, pp. 442–464, 2000.

**Hua Sun** (S'12-M'17) received his B.E. in Communications Engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 2011, M.S. in Electrical and Computer Engineering from University of California Irvine, USA, in 2013, and Ph.D. in Electrical Engineering from University of California Irvine, USA, in 2017. He is an Assistant Professor in the Department of Electrical Engineering at the University of North Texas, USA. His research interests include information theory and its applications to communications, privacy, networking, and storage.

Dr. Sun received the IEEE Jack Keil Wolf ISIT Student Paper Award in 2016, an IEEE GLOBECOM Best Paper Award in 2016, and the University of California Irvine CPCC Fellowship for the year 2011-2012. **Syed Ali Jafar** (S'99-M'04-SM'09-F'14) received his B. Tech. from IIT Delhi, India, in 1997, M.S. from Caltech, USA, in 1999, and Ph.D. from Stanford, USA, in 2003, all in Electrical Engineering. His industry experience includes positions at Lucent Bell Labs, Qualcomm Inc. and Hughes Software Systems. He is a Professor in the Department of Electrical Engineering and Computer Science at the University of California Irvine, Irvine, CA USA. His research interests include multiuser information theory, wireless communications and network coding.

Dr. Jafar is a recipient of the New York Academy of Sciences Blavatnik National Laureate in Physical Sciences and Engineering, the NSF CAREER Award, the ONR Young Investigator Award, the UCI Academic Senate Distinguished Mid-Career Faculty Award for Research, the School of Engineering Mid-Career Excellence in Research Award, the School of Engineering Maseeh Outstanding Research Award, the IEEE Information Theory Society Best Paper Award, IEEE Communications Society Best Tutorial Paper Award, IEEE Communications Society Heinrich Hertz Award, and three IEEE GLOBE-COM Best Paper Awards. His student co-authors received the IEEE Signal Processing Society Young Author Best Paper Award, and the Jack Wolf ISIT Best Student Paper Award. Dr. Jafar received the UC Irvine EECS Professor of the Year award six times, in 2006, 2009, 2011, 2012, 2014 and 2017 from the Engineering Students Council and the Teaching Excellence Award in 2012 from the School of Engineering. He was a University of Canterbury Erskine Fellow in 2010 and an IEEE Communications Society Distinguished Lecturer for 2013-2014. Dr. Jafar was recognized as a Thomson Reuters Highly Cited Researcher and included by Sciencewatch among The World's Most Influential Scientific Minds in 2014, 2015 and 2016. He served as Associate Editor for IEEE TRANSACTIONS ON COMMUNICATIONS 2004-2009, for IEEE COMMUNICATIONS LETTERS 2008-2009 and for IEEE TRANSACTIONS ON INFORMATION THEORY 2009-2012.