

# The Capacity of Single-Server Weakly-Private Information Retrieval

Hsuan-Yin Lin<sup>1</sup>, Senior Member, IEEE, Siddhartha Kumar<sup>2</sup>, Eirik Rosnes<sup>3</sup>, Senior Member, IEEE, Alexandre Graell i Amat<sup>4</sup>, Senior Member, IEEE, and Eitan Yaakobi<sup>5</sup>, Senior Member, IEEE

**Abstract**—A private information retrieval (PIR) protocol guarantees that a user can privately retrieve files stored in a database without revealing any information about the identity of the requested file. Existing information-theoretic PIR protocols ensure perfect privacy, i.e., zero information leakage to the servers storing the database, but at the cost of high download. In this work, we present weakly-private information retrieval (WPIR) schemes that trade off perfect privacy to improve the download cost when the database is stored on a single server. We study the tradeoff between the download cost and information leakage in terms of mutual information (MI) and maximal leakage (MaxL) privacy metrics. By relating the WPIR problem to rate-distortion theory, the download-leakage function, which is defined as the minimum required download cost of all single-server WPIR schemes for a given level of information leakage and a fixed file size, is introduced. By characterizing the download-leakage function for the MI and MaxL metrics, the capacity of single-server WPIR is fully described.

**Index Terms**—Private information retrieval, capacity, information-theoretic privacy, information leakage, single server.

## I. INTRODUCTION

USER privacy is becoming increasingly important both socially and politically in today's modern age of information (as demonstrated, for instance, by the European Union's General Data Protection Regulation). In this context, private information retrieval (PIR), introduced by Chor *et al.* [1], has gained traction in the information theory community. In PIR, a user can retrieve a file from a database without revealing the identity of the file to the servers storing it. From an information-theoretic perspective, the file size is typically much larger than the size of the queries to all servers.

Manuscript received July 27, 2020; revised December 4, 2020 and January 21, 2021; accepted January 24, 2021. Date of publication February 4, 2021; date of current version March 16, 2021. This work was supported in part by the Swedish Research Council under Grant 2016-04253; in part by the Israel Science Foundation under Grant 1817/18; and in part by the Technion Hiroshi Fujiwara Cyber Security Research Center and the Israel National Cyber Directorate. This article was presented in part at the IEEE International Symposium on Information Theory, Los Angeles, CA, USA, June 2020. (Corresponding author: Hsuan-Yin Lin.)

Hsuan-Yin Lin, Siddhartha Kumar, and Eirik Rosnes are with Simula UiB, 5006 Bergen, Norway (e-mail: lin@simula.no; kumarsi@simula.no; eirikrosnes@simula.no).

Alexandre Graell i Amat is with the Department of Electrical Engineering, Chalmers University of Technology, 41296 Gothenburg, Sweden, and also with Simula UiB, 5006 Bergen, Norway (e-mail: alexandre.graell@chalmers.se).

Eitan Yaakobi is with the Department of Computer Science, Technion—Israel Institute of Technology, Haifa 3200003, Israel (e-mail: yaakobi@gmail.com).

Digital Object Identifier 10.1109/JSAIT.2021.3056327

Therefore rather than accounting for both the upload and the download cost, as is usually done in the computer science community, here efficiency is measured in terms of the download cost. More precisely, efficiency is measured in terms of PIR rate, which is the ratio between the requested file size and the total number of symbols downloaded. The supremum of the PIR rate over all possible schemes and over all file sizes is called the PIR capacity.

PIR was first addressed in the information theory literature by Shah *et al.* [2], while the tradeoff between storage overhead and PIR rate was first considered in [3]. Shortly after, Sun and Jafar [4] characterized the PIR capacity for the classical PIR model of replicated servers. Since then the concept of PIR has been extended to several relevant scenarios: maximum distance separable (MDS) coded servers [5], [6], arbitrary linear coded servers [7]–[9], colluding servers [5], [7], [9]–[13], robust PIR [10], PIR with Byzantine servers [14], optimal upload cost of PIR, i.e., the smallest query information required to be sent to the servers [15], access complexity of PIR, i.e., the total number of symbols needed to be accessed across all servers in order to privately retrieve an arbitrary file [16], single-server PIR with side information [17], [18], PIR on graph-based replication systems [19], PIR with secure storage [20], [21], functional PIR codes [22], and private proximity retrieval codes [23].

Weakly-private information retrieval (WPIR) [24]–[26] is an interesting extension of the original PIR problem as it allows for improvements in the download cost at the expense of some information leakage to the servers on the identity of the requested file.<sup>1</sup> In particular, [24] considers the multi-server case with mutual information (MI) and worst-case information leakage [30] as privacy metrics, while [26] includes results for the maximum leakage (MaxL) privacy metric as well as converse bounds. In [25], Samy *et al.*, under the name of leaky PIR, consider a privacy metric related to *differential privacy* [31], [32] for the multi-server case. The multi-server case under the MaxL privacy metric has also been recently studied by Zhou *et al.* [33].

<sup>1</sup>On a related note, relaxing the requirement of perfect secrecy has been considered in the information theory literature in different contexts. For instance, in network coding, the term *weakly secure* is used when perfect security is only guaranteed for a subset of the messages multicasted from a source node [27], [28]. On the other hand, *weak security* in the context of secure communications refers to asymptotic *per-symbol* zero information leakage [29].

In the computer science literature, to the best of our knowledge, there are only a few works that have considered relaxing the perfect privacy requirement of PIR in order to improve performance. In the first work [34], which appeared almost two decades ago, the perfect privacy condition was relaxed by introducing the concept of *repudiation*. A protocol assures the repudiation property if the probability of all designed queries to retrieve any file stored in the database is strictly smaller than one. Hence, the user can deny any claim about the identity of the desired file by the server. However, the condition of repudiation can be achieved even if the server can determine the identity of the requested file almost surely, and thus, it does not provide a good level of information-theoretic privacy. More recently, Toledo *et al.* [35] adopted a privacy metric based on differential privacy in order to enhance the efficiency of PIR by lowering the level of privacy. Moreover, both works did not study any fundamental information-theoretic tradeoffs between information leakage and different costs under the considered privacy metrics. In contrast to WPIR, where the information leakage to the servers is considered, recently, the authors of [36] studied the information leakage of the nondesired files to the user in PIR systems. Furthermore, along the same lines of research, leaky PIR was extended to the symmetric PIR setting in [37]. Symmetric PIR is a variant of PIR where in addition to the privacy request, the user cannot learn anything about the remaining files in the database when the user retrieves its desired file [38]. In [37], the symmetric PIR requirement of zero information leakage on the nondesired files as well as the perfect privacy requirement on the identity of the requested file are both relaxed in order to improve the download rate.

In distributed storage systems (DSSs), for several applications, it is more realistic to assume that all servers can collude. In the single server scenario, it can be shown that all files stored in the database need to be downloaded to guarantee perfect privacy [39]. This implies that the PIR rate tends to zero as the number of stored files increases. In [17], the authors introduce PIR schemes that improve the download cost by leveraging on the assumption that the user has some prior side information on the content of the database. Two cases are considered, namely whether or not the privacy of the side information needs to be preserved. Lastly, latent-variable PIR for the single server setting was introduced in [40]. The goal of latent-variable PIR is to completely hide the latent attributes induced by the requested file identity. The latent-variable PIR framework can be seen as a variant of WPIR as no leakage on the latent attributes can be achieved even if parts of the identity of the requested file are leaked.

In this article, we relax the condition of perfect privacy in the single server setting. In similar lines to [24], we show that by relaxing the perfect privacy requirement, the download cost can be improved. Like [26], we consider both the MI and MaxL privacy metrics [41]–[43], where the latter is the most robust information-theoretic metric for information leakage yet known. In particular, we establish a connection between the single-server WPIR problem and rate-distortion theory, which

provides fundamental insights to describe the optimal tradeoff between the download cost and the allowed information leakage. The primary contribution of this work is to characterize the capacity (defined as the inverse of the minimum download cost over all possible schemes and over all file sizes) of single-server WPIR when the information leakage to the server is measured in terms of MI or MaxL. In this work, the minimum achievable download cost for a given information leakage constraint and for an *arbitrary* fixed file size is determined, and thus the WPIR capacity is derived. Especially, we propose a simple novel single-server WPIR scheme that achieves the WPIR capacity for both the MI and MaxL privacy metrics. Finally, we remark here that also the notion of differential privacy can be adapted to our setting, e.g., the local differential privacy metric [44], [45]. However, since the local differential privacy metric normally provides *stronger* privacy guarantees than the MI or MaxL metrics, it can readily be shown that it is not possible to further lower the download cost in the single server scenario. This is in contrast to the case of multiple servers [25].

The remainder of this article is structured as follows. Section II presents the notation, definitions, and the problem formulation. In Section III, we introduce the download-leakage function of single-server WPIR, which is defined as the minimum achievable download cost for a given information leakage constraint and for an arbitrary file size. Moreover, we discuss some properties of the function when the leakage is measured in terms of the MI or MaxL metrics. In Section IV-A, a basic solution for single-server WPIR is presented in which the file indices are partitioned into several partitions. In Section IV-B, we give a closed-form expression for the single-server WPIR capacity for both the MI and MaxL metrics. A capacity-achieving WPIR scheme is proposed in Section V. The converse result on the minimum download cost for the MI metric is provided in Section VI, while that of the MaxL metric is given in Section VII. Finally, Section VIII concludes this article.

## II. PRELIMINARIES AND PROBLEM STATEMENT

### A. Notation

We denote by  $\mathbb{N}$  the set of all positive integers,  $[a] \triangleq \{1, 2, \dots, a\}$ , and  $[a : b] \triangleq \{a, a + 1, \dots, b\}$  for  $a, b \in \{0\} \cup \mathbb{N}$  and  $a \leq b$ . The set of nonnegative real numbers is denoted by  $\mathbb{R}_+$ . Vectors are denoted by bold letters, matrices by sans serif capital letters, and sets by calligraphic uppercase letters, e.g.,  $\mathbf{x}$ ,  $\mathbf{X}$ , and  $\mathcal{X}$ , respectively. In general, vectors are represented as row vectors throughout this article. We use uppercase letters for random variables (RVs) (either scalar or vector), e.g.,  $X$  or  $\mathbf{X}$ . For a given index set  $\mathcal{S}$ , we write  $X^{\mathcal{S}}$  to represent  $\{X^{(m)} : m \in \mathcal{S}\}$ .  $X \perp\!\!\!\perp Y$  means that the two RVs  $X$  and  $Y$  are independent.  $(\cdot)^{\top}$  denotes the transpose of its argument. The Hamming weight of a vector  $\mathbf{x}$  is denoted by  $w_{\text{H}}(\mathbf{x})$ , while its support will be denoted by  $\chi(\mathbf{x})$ .  $\mathbb{E}_X[\cdot]$  and  $\mathbb{E}_{P_X}[\cdot]$  denote expectation with respect to the RV  $X$  and distribution  $P_X$ , respectively.  $\text{H}(X)$ ,  $\text{H}(P_X)$ , or  $\text{H}(p_1, \dots, p_{|\mathcal{X}|})$  represents the entropy of  $X$ , where  $P_X(\cdot) = (p_1, \dots, p_{|\mathcal{X}|})$  denotes the

distribution of the RV  $X$ .  $I(X; Y)$  denotes the MI between  $X$  and  $Y$ .

### B. System Model

We consider a single server that stores  $M$  independent files  $X^{(1)}, \dots, X^{(M)}$ , where each file  $X^{(m)} = (X_1^{(m)}, \dots, X_\beta^{(m)})$ ,  $m \in [M]$ , is represented as a length- $\beta$  row vector over  $\mathcal{X}$ . Assume that each element of  $X^{(m)}$  is chosen independently and uniformly at random from  $\mathcal{X}$ . Thus, we have  $H(X^{(m)}) = \beta \log_2 |\mathcal{X}|$  bits,  $\forall m \in [M]$ . A user wishes to efficiently retrieve  $X^{(M)}$  by allowing some information leakage to the server, where the requested file index  $M$  is assumed to be uniformly distributed over  $[M]$ .<sup>2</sup> Similar to the detailed mathematical description in [24], we give the following definition for a single-server WPIR scheme.

*Definition 1:* An  $M$ -file WPIR scheme  $\mathcal{C}$  for a single server storing  $M$  files consists of:

- A random strategy  $S$  with alphabet  $\mathcal{S}$ , which is privately designed by the user.
- A query function

$$\phi : \{1, \dots, M\} \times \mathcal{S} \rightarrow \mathcal{Q}$$

that generates a query  $\mathbf{Q} = \phi(M, S)$  with alphabet  $\mathcal{Q}$ , and induces a conditional probability mass function (PMF)  $P_{\mathbf{Q}|M}$ . The query  $\mathbf{Q}$  is sent to the server to retrieve the  $M$ -th file.

- An answer function

$$\varphi : \mathcal{Q} \times \mathcal{X}^{\beta M} \rightarrow \mathcal{A}^{\beta L}$$

that returns the answer  $A = \varphi(\mathbf{Q}, X^{[M]})$  back to the user, with download symbol alphabet  $\mathcal{A}$ . Here,  $L = L(\mathbf{Q})$  is the normalized length of the answer, which is a function of the query  $\mathbf{Q}$ .<sup>3</sup> More specifically, given a query realization  $\mathbf{Q} = \mathbf{q}$ ,  $L(\mathbf{q})$  can be seen as the codeword length of a code that encodes the files  $X^{[M]}$ , which is independent of the particular realization of the files.<sup>4</sup>

- A privacy leakage metric  $\rho^{(\cdot)}(P_{\mathbf{Q}|M}) \geq 0$ , which is defined as a function of  $P_{\mathbf{Q}|M}$ , that measures the amount of leaked information of the identity of the requested file to the server by observing the generated query  $\mathbf{Q}$ , where the superscript “ $(\cdot)$ ” indicates the used metric.

Furthermore, the scheme should allow a user to retrieve the requested file from the answer, the query, the index of the requested file, and the random strategy. In other words, this scheme must satisfy the condition of perfect retrievability,

$$H(X^{(M)} | A, \mathbf{Q}, M, S) = 0. \quad (1)$$

<sup>2</sup>Note that the requested file index  $M$  does not necessarily need to be uniformly distributed, which is referred to as semantic PIR in the literature [46].

<sup>3</sup>Note that in this work, the performance metric we focus on is the normalized download cost (see (4) later). Hence, without loss of generality, we define the answer-length in a normalized manner.

<sup>4</sup>From a source coding perspective, the files/sources are encoded by a *fixed-length* code, i.e.,  $L(\mathbf{q})$  is independent of the realization of the files, reflecting the fact that the files are independent and identically distributed (i.i.d.) according to a uniform distribution. As opposed to a *variable-length* code, for a fixed-length code all codewords are of equal length. Note that this setup follows the problem formulation in the PIR literature, see [15].

We remark that a PIR scheme corresponds to a WPIR scheme for which no information leakage is allowed.

### C. Metrics of Information Leakage

Given a single-server  $M$ -file WPIR scheme and a fixed distribution  $P_M$ , the conditional PMF of the query given the index  $M$  of the requested file,  $P_{\mathbf{Q}|M}$ , can be seen as a privacy mechanism (a randomized mapping). The server receives the random outcome  $\mathbf{Q}$  of the privacy mechanism  $P_{\mathbf{Q}|M}$ , and is curious about the index  $M$  of the requested file. The information leakage of a WPIR scheme is then measured with respect to its corresponding privacy mechanism  $P_{\mathbf{Q}|M}$ .

In this article, we focus on two commonly-used information-theoretic measures, namely MI and MaxL. For the former, the information leakage is quantified by

$$\rho^{(\text{MI})}(P_{\mathbf{Q}|M}) \triangleq I(M; \mathbf{Q}). \quad (2)$$

The second privacy metric, MaxL, which is introduced in [42], [43], is quantified by

$$\begin{aligned} \rho^{(\text{MaxL})}(P_{\mathbf{Q}|M}) &\triangleq \text{MaxL}(M; \mathbf{Q}) \\ &= \log_2 \sum_{q \in \mathcal{Q}} \max_{m \in [M]} P_{\mathbf{Q}|M}(q|m). \end{aligned} \quad (3)$$

We remark that MaxL can also be defined based on the *min-entropy* (*MinE*) information leakage  $I_\infty(M; \mathbf{Q})$  for the privacy mechanism  $P_{\mathbf{Q}|M}$ , where

$$I_\infty(M; \mathbf{Q}) \triangleq H_\infty(M) - H_\infty(M|\mathbf{Q})$$

and  $H_\infty(M)$  denotes the MinE measure that is widely discussed in the computer science literature, see [41]. On the other hand, the authors in [43] proposed the equivalent definition

$$\text{MaxL}(M; \mathbf{Q}) \triangleq \sup_{X \text{---} M \text{---} \mathbf{Q} \text{---} \hat{X}} \log_2 \frac{\Pr[X = \hat{X}]}{\max_{x \in \mathcal{X}} P_X(x)}$$

of MaxL, where the supremum is taken over all possible  $X$  and  $\hat{X}$  taking values in the same finite, but arbitrary alphabet  $\mathcal{X}$ , and the notation  $X \text{---} M \text{---} \mathbf{Q} \text{---} \hat{X}$  means that the RVs  $X$ ,  $M$ ,  $\mathbf{Q}$ , and  $\hat{X}$  form a Markov chain from left to right. Note that if  $\text{MaxL}(M; \mathbf{Q}) = \rho$  bits, the above definition indicates that for any possible randomized function  $X$  of  $M$ , the maximum probability of correctly guessing  $X$  based on  $\mathbf{Q}$  is bounded from above by the product of  $2^\rho$  and the maximum probability of guessing  $X$  with no observation.

It is worth mentioning that since we assume that  $M$  is uniformly distributed, the MinE information leakage and the MaxL privacy metric can be shown to be equivalent, i.e.,  $\text{MaxL}(M; \mathbf{Q}) = I_\infty(M; \mathbf{Q})$  [42, Th. 1], [43, Th. 1]. It is also worth mentioning that there is a relation between MaxL and differential privacy, see [42, Th. 3].

The following lemma summarizes some useful properties for both the MI and MaxL privacy metrics.

*Lemma 1 (Data Processing Inequalities [43, Lemma 1, Corollary 1]):* For any joint distribution  $P_{X,Y}$ ,

1) if the RVs  $X$ ,  $Y$ , and  $Z$  form a Markov chain, then

$$I(X; Z) \leq \min\{I(X; Y), I(Y; Z)\}, \text{ and} \\ \text{MaxL}(X; Z) \leq \min\{\text{MaxL}(X; Y), \text{MaxL}(Y; Z)\}.$$

2) Consider a fixed distribution  $P_X$ . Then, both  $I(X; Y)$  and  $\text{MaxL}(X; Y)$  are convex functions in  $P_{Y|X}$ .

Throughout this article, the information leakage metric of a WPIR scheme  $\mathcal{C}$  is denoted by  $\rho^{(\cdot)}(\mathcal{C})$ . Moreover, since  $\rho^{(\cdot)}$  is defined with a single argument of  $P_{Q|M}$ , we will also simply write the corresponding MI and MaxL metrics as a function of  $P_{Q|M}$ , i.e.,  $I(P_{Q|M}) = I(M; Q)$  and  $\text{MaxL}(P_{Q|M}) = \text{MaxL}(M; Q)$ .

#### D. Download Cost and Rate for a Single-Server WPIR Scheme

From the perfect privacy requirement of PIR, the server should not be able to differentiate the returned answers (e.g., by looking at their sizes) no matter which file index is requested. However, in contrast to PIR, the download cost for WPIR may be different for the retrieval of different files. Hence, the download cost is defined as the expected download cost over all possible requested files. The download cost of a single-server WPIR scheme  $\mathcal{C}$  for the retrieval of the  $m$ -th file, denoted by  $D^{(m)}(\mathcal{C})$ , is defined as the normalized expected length of the returned answer over all random queries,

$$D^{(m)}(\mathcal{C}) \triangleq \frac{\log_2 |\mathcal{A}| \mathbb{E}_{P_{Q|M=m}}[L(Q)]}{\log_2 |\mathcal{X}|},$$

and the overall download cost, denoted by  $D(\mathcal{C})$ , is measured in terms of the normalized expected download cost over all files, i.e.,

$$D(\mathcal{C}) \triangleq \frac{\log_2 |\mathcal{A}| \mathbb{E}_{P_M}[\mathbb{E}_{P_{Q|M}}[L(Q)]]}{\log_2 |\mathcal{X}|} \\ = \gamma \mathbb{E}_{M,Q}[L(Q)], \quad (4)$$

where  $\gamma \triangleq \log_2 |\mathcal{A}| / \log_2 |\mathcal{X}|$ . We remark that in general the sizes of the download symbol and file symbol alphabets can be different. However, for simplicity, we assume  $\gamma = 1$  throughout this article as the value of  $\gamma$  does not affect the generality of the results. Accordingly, the WPIR rate is defined as  $R(\mathcal{C}) \triangleq D(\mathcal{C})^{-1}$ .

Intuitively, a smaller download cost can be achieved if we allow a higher level of information leakage. In this article, our goal is to characterize the optimal tradeoff between the download cost and the allowed information leakage with respect to a privacy metric. We start with the following definition of an achievable download-leakage pair.

*Definition 2:* Consider a single server that stores  $M$  files. A download-leakage pair  $(D, \varrho)$  is said to be *achievable* in terms of the information leakage metric  $\rho^{(\cdot)}$  if there exists a WPIR scheme  $\mathcal{C}$  such that  $\mathbb{E}_{M,Q}[L(Q)] \leq D$  and  $\rho^{(\cdot)}(\mathcal{C}) \leq \varrho$ . The *download-leakage region* is the set of all achievable download-leakage pairs  $(D, \varrho)$ .

*Remark 1:* By Definition 2, it is clear that if the pair  $(D, \varrho)$  is achievable, then the pair  $(D', \varrho')$  with  $D' \geq D$  and  $\varrho' \geq \varrho$  is also achievable.

### III. CHARACTERIZATION OF THE OPTIMAL DOWNLOAD-LEAKAGE TRADEOFF

Consider the single-server WPIR problem with an arbitrary file size  $\beta$ , where the leakage is measured by  $\rho^{(\text{MI})}$  or  $\rho^{(\text{MaxL})}$ . The minimum achievable download cost for a given leakage constraint  $\varrho$  can be formulated as the optimization problem

$$\text{minimize } \mathbb{E}_{P_M P_{Q|M}}[L(Q)] \quad (5a)$$

$$\text{subject to } \{L(q)\}_{q \in \mathcal{Q}} \subseteq \mathcal{L}_{\text{ret}}, \quad (5b)$$

$$\rho^{(\cdot)}(P_{Q|M}) \leq \varrho \quad (5c)$$

over the query function  $\phi(\cdot)$  and the answer function  $\varphi(\cdot)$  from Definition 1, where  $\mathcal{L}_{\text{ret}}$  is defined as the set of the codeword lengths of all possible fixed-length codes that satisfy (1). In this way, it involves the query function, the answer generating function, and the decoding function from Definition 1. For brevity, the explicit query and answer function constraints are omitted as  $\mathcal{L}_{\text{ret}}$  implicitly involves these functions. We recall Definition 1 here that given a query realization  $Q = q$ ,  $L(q)$  can be seen as the codeword length of a fixed-length code that encodes the files and satisfies the *lossless* property in (1). Note again that the fixed-length assumption for the code, or equivalently that  $L(q)$  is independent of the specific realization of the files, reflects the fact that the files are i.i.d. according to a uniform distribution. Finally, we remark that since  $P_M$  is assumed to be fixed, the minimization over  $P_M P_{Q|M}$  in (5) is taken over the set of all conditional distributions  $P_{Q|M}$ .

#### A. Download-Leakage Function for Single-Server WPIR

To characterize the optimal achievable pairs of download cost and information leakage, we define two functions that describe the boundary of the download-leakage region.

*Definition 3:* For any file size  $\beta$ , the download-leakage function  $D^{(\cdot)}(\varrho)$  for single-server WPIR is the minimum of all possible download costs  $D$  for a given information leakage constraint  $\varrho$  such that  $(D, \varrho)$  is achievable, i.e.,

$$D^{(\cdot)}(\varrho) \triangleq \min_{\{L(q)\}_{q \in \mathcal{Q}} \subseteq \mathcal{L}_{\text{ret}}, P_{Q|M} : \rho^{(\cdot)}(P_{Q|M}) \leq \varrho} \mathbb{E}_{P_M P_{Q|M}}[L(Q)].$$

Following the notion of information-theoretic PIR capacity in the literature, we define the single-server WPIR capacity as the supremum of the inverse of all download-leakage functions over all possible values of  $\beta$  as follows.

*Definition 4:* The single-server WPIR capacity is defined as  $C^{(\cdot)}(\varrho) = \max_{\beta \in \mathbb{N}} \{[D^{(\cdot)}(\varrho)]^{-1}\}$ .

We remark here that the optimality results for the download-leakage functions in Sections VI and VII hold for any fixed file size  $\beta$ , since the solutions of the convex optimization problem formulations are independent of  $\beta$ , which indicates  $C^{(\cdot)}(\varrho) = [D^{(\cdot)}(\varrho)]^{-1}$  for any file size  $\beta$ . This also implies that increasing the file size does not further improve the performance. Note that it is known from the original work of Chor *et al.* [39] that the single-server PIR capacity is  $C_M = \frac{1}{M}$ .

We assume throughout the rest of this article that the perfect retrievability condition in (5b) holds. Hence, for convenience, we will sometimes drop the condition of (5b) in the download-leakage optimization formulation.

Naturally, we can also determine the optimal download-leakage region by interchanging the roles of the download cost and the information leakage.

*Definition 5:* For any file size  $\beta$ , the leakage-download function  $\rho^{(\cdot)}(\mathbf{D})$  for single-server WPIR is the minimum of all possible information leakages  $\varrho$  for a given download cost constraint  $\mathbf{D}$  such that  $(\mathbf{D}, \varrho)$  is achievable.

*Lemma 2:*

1) For any file size  $\beta$ , the MI download-leakage function

$$\mathbf{D}^{(\text{MI})}(\varrho) = \min_{P_{Q|M}: I(P_{Q|M}) \leq \varrho} \mathbb{E}_{P_M P_{Q|M}}[L(Q)]$$

is convex in  $\varrho$ .

2) For any file size  $\beta$ , the MaxL download-leakage function

$$\mathbf{D}^{(\text{MaxL})}(\varrho) = \min_{P_{Q|M}: \text{MaxL}(P_{Q|M}) \leq \varrho} \mathbb{E}_{P_M P_{Q|M}}[L(Q)]$$

is not a convex function, but  $\mathbf{D}^{(\text{MaxL})}(\log_2(\varrho))$  is convex in  $\varrho$ .

*Proof:* We first prove the lemma for MI leakage. Assume that  $\mathbf{D}^{(\text{MI})}(\varrho_1) = \mathbb{E}_{P_M P_{Q_1^*|M}}[L(Q)]$  and  $\mathbf{D}^{(\text{MI})}(\varrho_2) = \mathbb{E}_{P_M P_{Q_2^*|M}}[L(Q)]$  are achieved by the answer-lengths and conditional distributions  $(\{L(q)\}_{q \in \mathcal{Q}_1^*}, P_{Q_1^*|M})$  and  $(\{L(q)\}_{q \in \mathcal{Q}_2^*}, P_{Q_2^*|M})$ , respectively, where  $I(M; \mathcal{Q}_1^*) \leq \varrho_1$  and  $I(M; \mathcal{Q}_2^*) \leq \varrho_2$ . Let  $P_{Q_\lambda|M}$  be the distribution

$$P_{Q_\lambda|M} = (1 - \lambda)P_{Q_1^*|M} + \lambda P_{Q_2^*|M},$$

defined over  $\mathcal{Q}_\lambda \triangleq \{(i, q_i): q_i \in \mathcal{Q}_i, i = 1, 2\}$ . It can be seen that  $\{L(q)\}_{q \in \mathcal{Q}_\lambda} \subseteq \mathcal{L}_{\text{ret}}$ .

Observe that since  $I(M; \mathcal{Q})$  is convex in  $P_{Q|M}$ , it follows that  $I((1 - \lambda)P_{Q_1^*|M} + \lambda P_{Q_2^*|M}) \leq (1 - \lambda)I(P_{Q_1^*|M}) + \lambda I(P_{Q_2^*|M}) \leq (1 - \lambda)\varrho_1 + \lambda\varrho_2$ , which implies that  $P_{Q_\lambda|M}$  is an element of  $\{P_{Q|M}: I(P_{Q|M}) \leq (1 - \lambda)\varrho_1 + \lambda\varrho_2\}$ .

Thus, by definition we get

$$\begin{aligned} & \mathbf{D}^{(\text{MI})}((1 - \lambda)\varrho_1 + \lambda\varrho_2) \\ &= \min_{P_{Q|M}: I(P_{Q|M}) \leq (1 - \lambda)\varrho_1 + \lambda\varrho_2} \mathbb{E}_{P_M P_{Q|M}}[L(Q)] \\ &\leq \mathbb{E}_{P_M P_{Q_\lambda|M}}[L(Q)] \\ &\stackrel{(a)}{=} (1 - \lambda) \mathbb{E}_{P_M P_{Q_1^*|M}}[L(Q)] + \lambda \mathbb{E}_{P_M P_{Q_2^*|M}}[L(Q)] \\ &= (1 - \lambda)\mathbf{D}^{(\text{MI})}(\varrho_1) + \lambda\mathbf{D}^{(\text{MI})}(\varrho_2), \end{aligned}$$

where (a) follows directly from the definition of  $P_{Q_\lambda|M}$ . This shows that  $\mathbf{D}^{(\text{MI})}(\varrho)$  is convex in  $\varrho$  for the MI metric. The proof for the MaxL metric is analogous, since  $2^{\text{MaxL}(P_{Q|M})}$  is convex in  $P_{Q|M}$ . ■

From Remark 1 we can see that the convexity of the download-leakage function is very useful, since it can help to describe the download-leakage region if some achievable pairs are known. This observation can be summarized in the following corollary.

*Corollary 1:* Assume that both pairs  $(\mathbf{D}_1, \varrho_1)$  and  $(\mathbf{D}_2, \varrho_2)$  are achievable. Then, for any  $\lambda \in [0, 1]$ , the pair  $(\mathbf{D}_\lambda = (1 - \lambda)\mathbf{D}_1 + \lambda\mathbf{D}_2, \varrho_\lambda = (1 - \lambda)\varrho_1 + \lambda\varrho_2)$  is achievable under MI leakage, while the pair  $(\mathbf{D}_\lambda = (1 - \lambda)\mathbf{D}_1 + \lambda\mathbf{D}_2, \varrho_\lambda = \log_2[(1 - \lambda)2^{\varrho_1} + \lambda 2^{\varrho_2}])$  is achievable for MaxL.

## B. Connection to Rate-Distortion Theory

The celebrated *rate-distortion theory* of Shannon and Kolmogorov (see [47, Ch. 9], [48, Ch. 10], and references therein) determines the minimum source compression rate required to reproduce any source sequence under a fidelity constraint, which is provided through a *distortion measure* between the source sequence and the reconstructed sequence. Consider an information source sequence with i.i.d. components according to  $P_X$  and a distortion measure  $d(\mathbf{x}, \hat{\mathbf{x}})$  between the source sequence  $\mathbf{x}$  and the reconstructed sequence  $\hat{\mathbf{x}}$ . The optimal rate-distortion region is characterized by the *rate-distortion function*, defined as the minimum achievable compression rate  $I(X; \hat{X})$  under a given constraint on the average distortion  $\mathbb{E}_{P_X P_{\hat{X}|X}}[d(X, \hat{X})]$ , where  $\hat{X}$  represents the reconstructed source.

One important observation from (4) is that, if we add the desired file index  $m$  as an argument to the answer-length function  $L$  by defining  $L(m, \mathbf{q}) \triangleq L(\mathbf{q})$  for all  $m \in [M]$  for which  $P_{Q|M}(\mathbf{q}|m) > 0$ , and  $L(m, \mathbf{q}) \triangleq \infty$  otherwise (i.e., an infinite length for a given  $m$  and query realization  $\mathbf{q}$  indicates that  $\mathbf{q}$  is never sent when requesting the  $m$ -th file), then the download cost can be expressed as  $\mathbb{E}_{P_M P_{Q|M}}[L(Q)] = \mathbb{E}_{P_M P_{Q|M}}[L(M, Q)]$ . Thus, in terms of the MI privacy metric, the leakage-download function of a given WPIR scheme can be related to the rate-distortion function, where the leakage and the download cost play similar roles as the compression rate and the average distortion, respectively. Below, we will equivalently use either  $L(\mathbf{q})$  or  $L(m, \mathbf{q})$  (as defined above). We defer the detailed discussion to Section VI where we will utilize results from rate-distortion theory to characterize the optimal leakage-download tradeoff for single-server WPIR.

## IV. RESULTS

### A. Partition WPIR Scheme

In [24], a WPIR scheme based on partitioning was proposed. The set of all file indices, i.e.,  $[M]$ , is first privately pre-partitioned into  $\eta$  equally-sized partitions by the user, each consisting of  $M_\eta$  file indices, where  $M_\eta = M/\eta \in \mathbb{N}$ . If there exists a viable  $M_\eta$ -file WPIR scheme, the user can apply the  $M_\eta$ -file WPIR scheme as a subscheme on each partition, and retrieve a file from the corresponding partition.

The partition  $M$ -file WPIR scheme is formally described as follows. Assume that the requested file  $\mathbf{X}^{(m)}$  belongs to the  $j$ -th partition, where  $j \in [\eta]$ . Then, the query  $\mathcal{Q}$  is constructed as

$$\mathcal{Q} = (\tilde{\mathcal{Q}}, j) \in \tilde{\mathcal{Q}} \times [\eta], \quad (6)$$

where  $\tilde{\mathcal{Q}}$  is the query of an existing  $M_\eta$ -file WPIR scheme.

The following theorem states the achievable download-leakage pairs of the partition scheme.

*Theorem 1:* Consider a single server that stores  $M$  files and let  $M_\eta = M/\eta \in \mathbb{N}$ ,  $\eta \in \mathbb{N}$ . Assume that an  $M_\eta$ -file WPIR scheme  $\tilde{\mathcal{C}}$  with achievable download-leakage pair  $(\tilde{\mathbf{D}}, \tilde{\varrho})$  exists. Then, the download-leakage pair

$$(\mathbf{D}(\mathcal{C}), \rho^{(\cdot)}(\mathcal{C})) = (\tilde{\mathbf{D}}, \tilde{\varrho} + \log_2 \eta) \quad (7)$$

is achievable by the  $M$ -file partition scheme  $\mathcal{C}$  constructed from  $\tilde{\mathcal{C}}$  as described in (6).

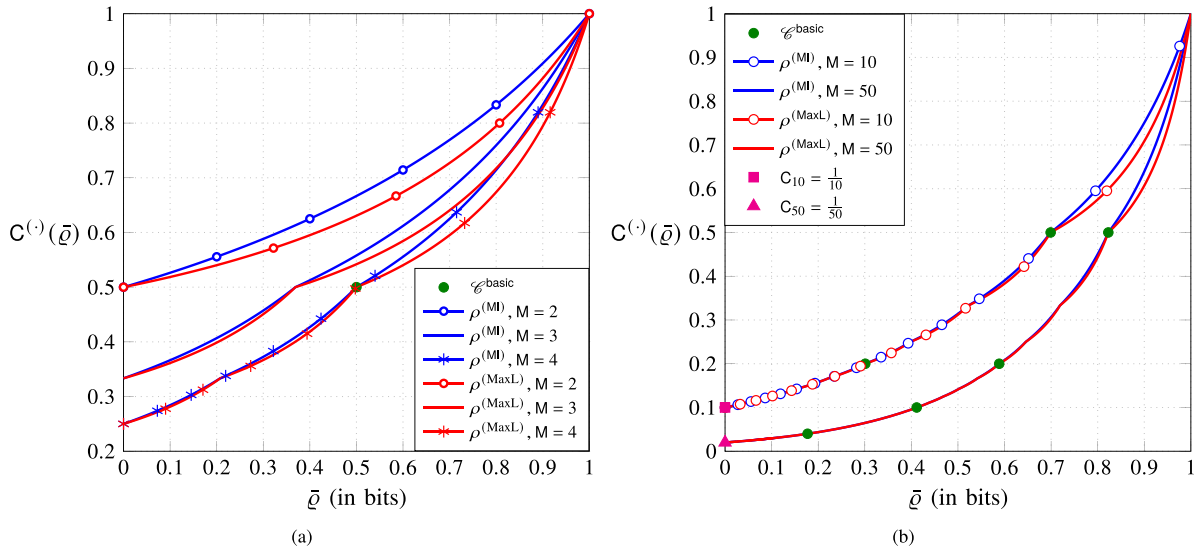


Fig. 1. (a) The capacity  $C^{(\cdot)}(\bar{q})$  for a small number of files  $M = 2, 3, 4$  with privacy metrics  $\rho^{(MI)}$  and  $\rho^{(MaxL)}$ . (b) The capacity  $C^{(\cdot)}(\bar{q})$  for a large number of files  $M = 10, 50$  with privacy metrics  $\rho^{(MI)}$  and  $\rho^{(MaxL)}$ . The dark green circles mark the achievable rate-leakage pairs of  $\mathcal{L}^{basic}$ .

*Proof:* The theorem for the MI privacy metric is proved in [26] (see proof of Theorem 2). Here, we provide the proof for the MaxL metric, which follows a similar argumentation as in the proof for the MI metric.

We refer to the requested file index  $M$  by the pair  $(\tilde{M}, j)$ , where  $\tilde{M}$  represents the requested file index in the  $j$ -th partition,  $j \in [\eta]$ . Hence, from (3), we have

$$\begin{aligned} 2^{\text{MaxL}(M; \mathcal{Q})} &= \sum_{q \in \mathcal{Q}} \max_{m \in [M]} P_{\mathcal{Q}|M}(q|m) \\ &= \sum_{j \in [\eta]} \sum_{\tilde{q} \in \tilde{\mathcal{Q}}} \max_{m \in [M]} P_{\mathcal{Q}|M}(q|m) \\ &\stackrel{(a)}{=} \sum_{j \in [\eta]} \sum_{\tilde{q} \in \tilde{\mathcal{Q}}} \max_{m \in [M_\eta]} P_{\tilde{\mathcal{Q}}|\tilde{M}}(\tilde{q}|m) \\ &= \eta \cdot 2^{\text{MaxL}(\tilde{M}; \tilde{\mathcal{Q}})} \leq 2^{\log_2 \eta + \bar{q}}, \end{aligned}$$

where (a) follows since for the  $j$ -th partition, the conditional PMF  $P_{\mathcal{Q}|M}$  is equal to  $P_{\tilde{\mathcal{Q}}|\tilde{M}}$  of the  $M_\eta$ -file WPIR scheme. Using a similar argumentation as above, it can also be verified that  $D(\mathcal{C}) = D(\tilde{\mathcal{C}}) \leq \bar{D}$ . ■

Since a PIR scheme is also a WPIR scheme, this simple approach for the construction of WPIR schemes can also be adapted to use any of the existing  $M_\eta$ -file PIR schemes in the literature as a subscheme. We refer to the partition scheme that uses a PIR scheme as the underlying subscheme and the query generation in (6) as a *basic scheme* and denote it by  $\mathcal{L}^{basic}$  (it achieves the pair in (7) with  $\bar{q} = 0$ ). It can be seen that for the single-server setting, the basic scheme simply retrieves all the files in the partition that includes the requested file. This idea will be extended to our capacity-achieving scheme presented in Section V, where for any subset  $\mathcal{M} \subseteq [M]$  that includes the requested file, all files in  $\mathcal{M}$  are downloaded.

### B. The Capacity of Single-Server WPIR

The main result of this work is the characterization of the optimal tradeoff between the download cost and the information leakage for single-server WPIR for an arbitrary

number of files and a fixed file size  $\beta$  for the MI and MaxL privacy metrics. The capacity of single-server WPIR for the MI privacy metric is stated in the following theorem. For the sake of illustration, we consider the normalized leakage metric  $\bar{\rho}^{(\cdot)} \triangleq \frac{\rho^{(\cdot)}}{\log_2 M}$ .

**Theorem 2:** For a single server that stores  $M$  files, the WPIR capacity for the MI leakage metric  $\rho^{(MI)}$  is

$$\begin{aligned} C^{(MI)}(\bar{q}) &= \left[ w + \frac{\log_2 \frac{M}{w}}{\log_2 \frac{w}{w-1}} - \frac{\bar{q} \log_2 M}{\log_2 \frac{w}{w-1}} \right]^{-1}, \\ \text{for } 1 - \frac{\log_2 w}{\log_2 M} &\leq \bar{q} \leq 1 - \frac{\log_2 (w-1)}{\log_2 M}, \quad w \in [2 : M]. \quad (8) \end{aligned}$$

The following theorem states the single-server WPIR capacity for the MaxL privacy metric.

**Theorem 3:** For a single server that stores  $M$  files, the WPIR capacity for the MaxL metric  $\rho^{(MaxL)}$  is

$$\begin{aligned} C^{(MaxL)}(\bar{q}) &= \left[ w + \frac{\frac{M}{w}}{\frac{M}{w-1} - \frac{M}{w}} - \frac{2\bar{q} \log_2 M}{\frac{M}{w-1} - \frac{M}{w}} \right]^{-1}, \\ \text{for } 1 - \frac{\log_2 w}{\log_2 M} &\leq \bar{q} \leq 1 - \frac{\log_2 (w-1)}{\log_2 M}, \quad w \in [2 : M]. \end{aligned}$$

The achievability proof of Theorems 2 and 3 appears in Section V while the converse part appears in Sections VI and VII, for Theorems 2 and 3, respectively.

For the MI and MaxL privacy metrics, the achievable rate-leakage pairs of  $\mathcal{L}^{basic}$  in (7) and the capacity  $C^{(\cdot)}(\bar{q})$  for different number of files  $M$ , are depicted in Fig. 1.

It is worthwhile noting that the curve  $[C^{(\cdot)}(\cdot)]^{-1}$  is a piecewise continuous function. For the MI privacy metric, (8) indicates that the minimum download cost  $[C^{(MI)}(\cdot)]^{-1}$  is a piecewise linear function in  $\bar{\rho}$ , illustrating the convexity of Lemma 2. Note also that, when  $M_\eta = M/\eta \in \mathbb{N}$ , the basic scheme  $\mathcal{L}^{basic}$  achieves the capacity for both the MI and MaxL privacy metrics.

Next, we consider the asymptotic capacity of single-server WPIR, i.e., the capacity as the number of files  $M$  tends to

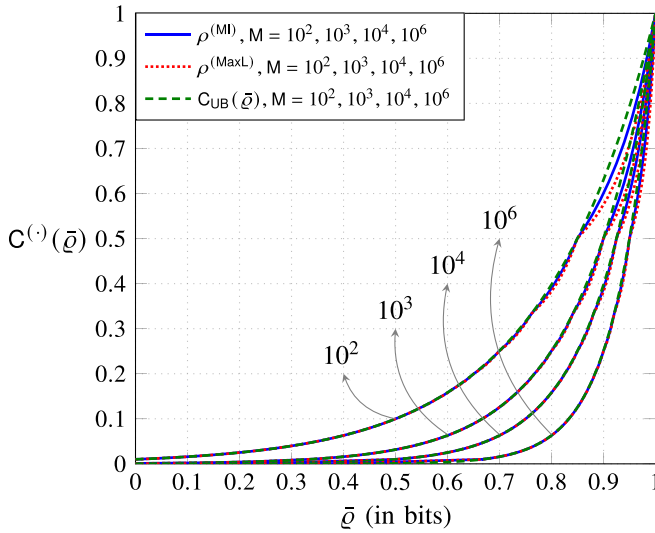


Fig. 2. The capacity  $C^{(\cdot)}(\bar{\rho})$  and its upper bound  $C_{UB}(\bar{\rho})$  for a number of files  $M = 10^2, 10^3, 10^4,$  and  $10^6$  with privacy metrics  $\rho^{(MI)}$  and  $\rho^{(MaxL)}$ .

infinity. An upper bound on the single-server WPIR capacity for any number of files is given in the following theorem.

*Theorem 4:* For a single server that stores  $M$  files, the single-server WPIR capacity under both the MI metric  $\rho^{(MI)}$  and the MaxL metric  $\rho^{(MaxL)}$  is bounded from above by

$$C^{(\cdot)}(\bar{\rho}) \leq C_{UB}(\bar{\rho}) \triangleq \frac{1}{M^{1-\bar{\rho}}}, \quad 0 \leq \bar{\rho} \leq 1.$$

*Proof:* Since it is easy to show that the capacity for MI leakage is larger than or equal to the capacity for MaxL (see Fig. 1), we only need to prove that the inverse of (8) is bounded from below by  $M^{1-\bar{\rho}}$ . Observe that  $M^{1-\bar{\rho}}$  is a convex function of  $\bar{\rho}$  and each point  $(w, 1 - \frac{\log_2 w}{\log_2 M})$ ,  $w \in [M]$ , lies on the curve described by that function, i.e.,  $M^{1-\bar{\rho}} = w$  for  $\bar{\rho} = 1 - \frac{\log_2 w}{\log_2 M}$ . Thus, by the convexity of  $M^{1-\bar{\rho}}$ , we have  $M^{1-\bar{\rho}} \leq [C^{(MI)}(\bar{\rho})]^{-1}$ , where the inequality follows since the inverse of (8) can be seen as a convex combination of  $1 - \frac{\log_2 w}{\log_2 M}$  and  $1 - \frac{\log_2(w-1)}{\log_2 M}$ ,  $w \in [2 : M]$ . ■

In Fig. 2, the capacity  $C^{(\cdot)}(\bar{\rho})$  and the upper bound  $C_{UB}(\bar{\rho})$  are plotted for  $M = 10^2, 10^3, 10^4,$  and  $10^6$ , which illustrates the asymptotic behavior of  $C^{(\cdot)}(\bar{\rho})$  as  $M$  tends to infinity. Note that since one can simply download the requested file from the server in the special case of  $\bar{\rho} = 1$ , and the WPIR rate must be smaller than or equal to 1, we have  $C^{(\cdot)}(1) = 1$  for either a finite or infinite number of files  $M$ . Hence, by Theorem 4 it can be shown that as  $M$  tends to infinity, the asymptotic capacity is equal to

$$C_{\infty}^{(\cdot)}(\bar{\rho}) = \begin{cases} 0 & \text{if } 0 \leq \bar{\rho} < 1, \\ 1 & \text{if } \bar{\rho} = 1. \end{cases}$$

This indicates that the asymptotic capacity is still equal to zero, unless the server exactly knows the index of the requested file.

## V. ACHIEVABILITY

Throughout this section, for simplicity, we set  $\beta = 1$ . Hence, from Section II-B, we have  $X^{(m)} = X_1^{(m)}$  and  $H(X^{(m)}) = \log_2 |\mathcal{X}|$  bits,  $\forall m \in [M]$ . In fact, our proposed

TABLE I  
CONDITIONAL PMFS  $P_{Q_w|M}$ ,  $w \in [3]$

$Q_1$	$P_{Q_1 M}(q 1)$	$P_{Q_1 M}(q 2)$	$P_{Q_1 M}(q 3)$	$A_1$	$P_{Q_1}(q)$
(1, 0, 0)	1	0	0	$X_1^{(1)}$	$\frac{1}{3}$
(0, 1, 0)	0	1	0	$X_1^{(2)}$	$\frac{1}{3}$
(0, 0, 1)	0	0	1	$X_1^{(3)}$	$\frac{1}{3}$

$Q_2$	$P_{Q_2 M}(q 1)$	$P_{Q_2 M}(q 2)$	$P_{Q_2 M}(q 3)$	$A_2$	$P_{Q_2}(q)$
(1, 1, 0)	$\frac{1}{2}$	$\frac{1}{2}$	0	$\{X_1^{(1)}, X_1^{(2)}\}$	$\frac{1}{3}$
(1, 0, 1)	$\frac{1}{2}$	0	$\frac{1}{2}$	$\{X_1^{(1)}, X_1^{(3)}\}$	$\frac{1}{3}$
(0, 1, 1)	0	$\frac{1}{2}$	$\frac{1}{2}$	$\{X_1^{(2)}, X_1^{(3)}\}$	$\frac{1}{3}$

$Q_3$	$P_{Q_3 M}(q 1)$	$P_{Q_3 M}(q 2)$	$P_{Q_3 M}(q 3)$	$A_3$	$P_{Q_3}(q)$
(1, 1, 1)	1	1	1	$\{X_1^{(1)}, X_1^{(2)}, X_1^{(3)}\}$	1

single-server WPIR capacity-achieving scheme can be easily generalized to an arbitrary file size  $\beta$ , which indicates that subpacketization does not improve the performance of single-server WPIR. We will later show that this scheme is optimal for both the MI and MaxL privacy metrics.

### A. Motivating Example: $M = 3$ Files

Before describing the achievable scheme in detail for the general case of  $M$  files, we present an example for  $M = 3$ . Assume that the single server stores  $M = 3$  files,  $X_1^{(1)}$ ,  $X_1^{(2)}$ , and  $X_1^{(3)}$ . We design the queries and answers via a conditional distribution  $P_{Q_w|M}$ ,  $w \in [3]$ , defined in Table I. It can be easily verified that the perfect retrievability condition of (1) is satisfied for the three PMFs. Moreover, the download-leakage pairs  $(D_w, Q_w) = (w, \log_2 \frac{3}{w})$  are achievable, by  $P_{Q_w|M}$ ,  $w \in [3]$ , in terms of the MI or MaxL privacy metrics.

Now, construct two conditional query distributions as follows,

$$P_{Q_{\lambda_1}|M} = (1 - \lambda_1)P_{Q_2|M} + \lambda_1 P_{Q_1|M}, \quad (9)$$

$$P_{Q_{\lambda_2}|M} = (1 - \lambda_2)P_{Q_3|M} + \lambda_2 P_{Q_2|M}, \quad (10)$$

where  $0 \leq \lambda_1, \lambda_2 \leq 1$ .

The retrievability condition can be easily verified for the WPIR scheme defined by (9)–(10). Using (4), (2), and (3), respectively, and the conditional PMFs listed in Table I (or, alternatively, Corollary 1), it follows that the WPIR scheme defined by (9)–(10) achieves the download cost

$$D(\mathcal{C}) = \begin{cases} (1 - \lambda_1)D_2 + \lambda_1 D_1 = 2 - \lambda_1, & 0 \leq \lambda_1 \leq 1, \\ (1 - \lambda_2)D_3 + \lambda_2 D_2 = 3 - \lambda_2, & 0 \leq \lambda_2 \leq 1, \end{cases}$$

the MI leakage

$$\rho^{(MI)} = \begin{cases} (1 - \lambda_1) \log_2 \frac{3}{2} + \lambda_1 \log_2 \frac{3}{1}, & 0 \leq \lambda_1 \leq 1, \\ (1 - \lambda_2) \log_2 \frac{3}{3} + \lambda_2 \log_2 \frac{3}{2}, & 0 \leq \lambda_2 \leq 1, \end{cases}$$

and the MaxL

$$\rho^{(MaxL)} = \begin{cases} \log_2 \left( (1 - \lambda_1) \frac{3}{2} + \lambda_1 \frac{3}{1} \right), & 0 \leq \lambda_1 \leq 1, \\ \log_2 \left( (1 - \lambda_2) \frac{3}{3} + \lambda_2 \frac{3}{2} \right), & 0 \leq \lambda_2 \leq 1. \end{cases}$$

In terms of the MI or MaxL privacy metrics, it can be verified that the download cost corresponds to the single-server WPIR capacity for  $M = 3$ . Note that for  $M > 2$ , the capacity is a piecewise continuous function (see Fig. 1(a)).

### B. Arbitrary Number of Files $M$

We describe the achievable scheme for the general case of  $M$  files. From Corollary 1, it follows that it is sufficient to show that the download-leakage pairs

$$(\mathbf{D}_w, \mathcal{Q}_w) = \left( w, \log_2 \frac{M}{w} \right), \quad w \in [M],$$

are achievable.

1) *Query Generation*: Consider  $M$  random queries  $\mathcal{Q}_w$ ,  $w \in [M]$ , whose alphabet is  $\mathcal{Q}_w \triangleq \{\mathbf{q} = (q_1, \dots, q_M) \in \{0, 1\}^M : w_H(\mathbf{q}) = w\}$ . Recall here that  $\chi(\mathbf{q})$  denotes the support of a vector  $\mathbf{q}$ . Given any requested file index  $m \in [M]$ , each query  $\mathbf{q} \in \mathcal{Q}_w$  sent to the server is generated by the conditional PMF

$$P_{\mathcal{Q}_w|M}(\mathbf{q}|m) = \begin{cases} \frac{1}{\binom{M-1}{w-1}} & \text{if } |\chi(\mathbf{q}) \setminus \{m\}| = w - 1, \\ 0 & \text{otherwise.} \end{cases}$$

This is clearly a valid query design, since for each  $m \in [M]$ , we have  $\sum_{\mathbf{q} \in \mathcal{Q}_w} P_{\mathcal{Q}_w|M}(\mathbf{q}|m) = 1$ .

2) *Answer Construction*: The answer function  $\varphi$  maps the query  $\mathbf{q} \in \mathcal{Q}_w$  onto  $\mathbf{A} = \varphi(\mathbf{q}, \mathbf{X}^{[M]}) = X_1^{\chi(\mathbf{q})}$ . The answer length is  $L(\mathbf{q}) = w$ .

3) *Download Cost and Information Leakage*: Clearly, the download cost is equal to  $\mathbb{E}_{P_M P_{\mathcal{Q}_w|M}}[L(\mathcal{Q}_w)] = w$ . The MI leakage is

$$\begin{aligned} \rho^{(MI)}(P_{\mathcal{Q}_w|M}) &= I(M; \mathcal{Q}_w) = H(M) - H(M|\mathcal{Q}_w) \\ &= \log_2 M - \log_2 w = \log_2 \frac{M}{w} \end{aligned}$$

and the MaxL is

$$\begin{aligned} \rho^{(MaxL)}(P_{\mathcal{Q}_w|M}) &= \log_2 \sum_{\mathbf{q} \in \mathcal{Q}_w} \max_{m \in [M]} \frac{1}{\binom{M-1}{w-1}} \\ &= \log_2 \frac{\binom{M}{w}}{\binom{M-1}{w-1}} = \log_2 \frac{M}{w}. \end{aligned}$$

This completes the achievability proof of Theorems 2 and 3.

Notice that the presented capacity-achieving WPIR scheme can be seen as a generalization of the basic WPIR scheme  $\mathcal{C}^{\text{basic}}$ . If  $M/\eta = w \in \mathbb{N}$ , since  $(\tilde{\mathbf{D}}, \tilde{\varrho}) = (w, 0)$  is achievable for a  $w$ -file single-server PIR scheme, from Theorem 1 it follows that the download-leakage pair  $(\mathbf{D}(\mathcal{C}^{\text{basic}}), \rho^{(\cdot)}(\mathcal{C}^{\text{basic}})) = (w, \log_2 \eta) = (w, \log_2 \frac{M}{w})$  is also achievable by  $\mathcal{C}^{\text{basic}}$  for both the MI and MaxL privacy metrics.

## VI. CONVERSE OF THEOREM 2

For any file size  $\beta$ , a general converse (upper bound) can be derived from the download-leakage function of a given leakage constraint  $\varrho$ , or equivalently, from the leakage-download function of a given download cost constraint  $\mathbf{D}$ . The proof consists

of two parts, and we start by outlining the main arguments of each part before diving further into the technical details.

*Part 1*: Consider an arbitrary WPIR scheme. Without loss of optimality, for any query sent to the server, the answer from the server can be assumed to be a subset of the files that includes the desired file. This is because in order to have perfect retrievability, downloading any linear combinations, or any coded forms of a subset of the files can only lead to a higher download cost. Moreover, it can also not increase the privacy leakage, because the server can only infer the identity of the desired file from an answer that is able to recover a subset of the files. We formally prove this argument later in this section. Note that this part is also used in the converse proof of Theorem 3 as it holds for both the MI and MaxL privacy metrics.

*Part 2*: From Part 1 we can limit the consideration to schemes for which all answers are subsets of files that include the desired file. Since the minimum achievable information leakage for a given download cost constraint among this limited family of schemes can be related to the rate-distortion function with a certain distortion measure, we can apply a known lower bound on the rate-distortion function in order to find an optimal scheme from this family. Finally, we show that the optimal scheme is exactly the scheme we propose in Section V. Thus, for a given download cost constraint, the leakage of any WPIR scheme is bounded below by the leakage of the scheme proposed in Section V.

We start to prove the first part of the converse proof. Given an arbitrary query set  $\mathcal{Q}$  of a WPIR scheme with  $\{L(\mathbf{q})\}_{\mathbf{q} \in \mathcal{Q}} \subseteq \mathcal{L}_{\text{ret}}$ , similar to (5), the minimum leakage of this WPIR scheme for a given download cost constraint  $\mathbf{D}$  can be formulated as the convex optimization problem

$$\text{minimize } I(M; \mathcal{Q}) \quad (11a)$$

$$\text{subject to } \mathbb{E}_{P_M P_{\mathcal{Q}|M}}[L(\mathcal{Q})] \leq \mathbf{D}. \quad (11b)$$

The minimization is taken over the set of all conditional distributions  $P_{\mathcal{Q}|M}$  such that (11b) is satisfied, namely the set

$$\mathcal{F}_{\mathbf{D}} = \left\{ P_{\mathcal{Q}|M} : \sum_{\mathbf{q}} \sum_{m \in [M]} P_M(m) P_{\mathcal{Q}|M}(\mathbf{q}|m) L(m, \mathbf{q}) \leq \mathbf{D} \right\}.$$

Next, we show a lower bound to (11) by defining a new RV that is a function of  $\mathcal{Q}$ . To facilitate the exposition, we introduce the following notation.

- For any nonempty subset  $\mathcal{M} \subseteq [M]$ , we define  $\tilde{\mathcal{Q}}^{\mathcal{M}}$  to be the set of queries that are designed to recover the files  $\mathbf{X}^{(m)}$ ,  $m \in \mathcal{M}$ , i.e.,  $\tilde{\mathcal{Q}}^{\mathcal{M}} \triangleq \{\mathbf{q} \in \mathcal{Q} : H(\mathbf{X}^{(m)} | \mathbf{A}, \mathcal{Q} = \mathbf{q}, M = m, S = s) = 0, \forall m \in \mathcal{M}\}$ . Furthermore, define

$$\mathcal{Q}^{\mathcal{M}} \triangleq \tilde{\mathcal{Q}}^{\mathcal{M}} \setminus \left( \bigcup_{\mathcal{M}' \subseteq [M] \setminus \mathcal{M}} \tilde{\mathcal{Q}}^{\mathcal{M}'} \right).$$

The set  $\mathcal{Q}^{\mathcal{M}}$  contains all queries that are designed to recover all files in  $\mathcal{M}$ , but no more. Note that if  $\mathbf{q} \in \mathcal{Q}^{\mathcal{M}}$  but  $m \notin \mathcal{M}$ , then  $P_{\mathcal{Q}|M}(\mathbf{q}|m) = 0$ .



- A binary length- $M$  *indicator vector*  $\mathbf{1}_{\mathcal{M}} = (u_1, \dots, u_M)$  of a subset  $\mathcal{M} \subseteq [M]$ ,  $\mathcal{M} \neq \emptyset$ , is defined as

$$u_m = \begin{cases} 1 & \text{if } m \in \mathcal{M}, \\ 0 & \text{otherwise.} \end{cases}$$

$\mathbf{Q}$  is a RV that is induced by the conditional distribution  $P_{\mathbf{Q}|M}$ . Further define a new RV  $\mathbf{U} = f(\mathbf{Q})$ , where

$$f(\mathbf{q}) \triangleq \mathbf{1}_{\mathcal{M}}, \quad \text{for } \mathbf{q} \in \mathcal{Q}^M, \forall \mathcal{M} \subseteq [M], \mathcal{M} \neq \emptyset. \quad (12)$$

Note that the mapping in (12) is well-defined since the query sets  $\mathcal{Q}^M$  are disjoint and their union is equal  $\mathcal{Q}$ , i.e., they constitute a partition of  $\mathcal{Q}$ . This leads to the conditional PMF

$$P_{U|M}(\mathbf{u}|m) = \begin{cases} \sum_{\mathbf{q} \in \mathcal{Q}^{\chi(\mathbf{u})}} P_{\mathbf{Q}|M}(\mathbf{q}|m) & \text{if } m \in \chi(\mathbf{u}), \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

To further simplify the notation, in the following we use  $p(\mathbf{q}|m)$  and  $p(\mathbf{u}|m)$  to denote the conditional PMFs  $P_{\mathbf{Q}|M}$  and  $P_{U|M}$ , respectively.

We then design a WPIR scheme where the queries are generated according to  $p(\mathbf{u}|m)$  and the normalized answer-length function is constructed for any  $\mathbf{u} \in \{0, 1\}^M$  as

$$L(m, \mathbf{u}) = \begin{cases} w_{\mathbf{H}}(\mathbf{u}) & \text{if } m \in \chi(\mathbf{u}), \\ \infty & \text{otherwise,} \end{cases}$$

where an infinite length for a given file index  $m$  and query  $\mathbf{u}$  indicates that the  $m$ -th file is never retrieved by the designed query  $\mathbf{u}$ .

Accordingly, we define the set

$$\mathcal{P}_{\mathbf{D}} = \left\{ p(\mathbf{q}|m) : \sum_{\mathbf{u} \in \{0, 1\}^M} \sum_{\mathbf{q} \in \mathcal{Q}^{\chi(\mathbf{u})}} \sum_{m \in [M]} P_M(m) p(\mathbf{q}|m) L(m, \mathbf{u}) \leq \mathbf{D} \right\}.$$

Given an arbitrary  $p(\mathbf{q}|m) \in \mathcal{F}_{\mathbf{D}}$ ,

$$\begin{aligned} & \sum_{\mathbf{u} \in \{0, 1\}^M} \sum_{m \in [M]} P_M(m) p(\mathbf{u}|m) L(m, \mathbf{u}) \\ &= \sum_{\mathbf{u} \in \{0, 1\}^M} \sum_{m \in [M]} P_M(m) \sum_{\mathbf{q} \in \mathcal{Q}^{\chi(\mathbf{u})}} p(\mathbf{q}|m) L(m, \mathbf{u}) \\ &\stackrel{(a)}{\leq} \sum_{\mathbf{u} \in \{0, 1\}^M} \sum_{\mathbf{q} \in \mathcal{Q}^{\chi(\mathbf{u})}} \sum_{m \in [M]} P_M(m) p(\mathbf{q}|m) L(m, \mathbf{q}) \\ &= \sum_{\mathbf{q} \in \mathcal{Q}} \sum_{m \in [M]} P_M(m) p(\mathbf{q}|m) L(m, \mathbf{q}) \leq \mathbf{D}, \end{aligned}$$

where (a) holds because for any  $\mathbf{q} \in \mathcal{Q}^{\chi(\mathbf{u})}$  we know from the lossless source coding theorem that the normalized length of the answer  $L(m, \mathbf{q})$  is larger than or equal to the ratio between the total sizes of the retrieved  $w_{\mathbf{H}}(\mathbf{u}) = L(m, \mathbf{u})$  files and the logarithm of the source code's alphabet, i.e., it satisfies  $L(m, \mathbf{q}) \geq L(m, \mathbf{u}) \log_2 |\mathcal{X}| / \log_2 |\mathcal{A}| = L(m, \mathbf{u})$ .<sup>5</sup> Hence, it follows

<sup>5</sup>The fundamental theorem of lossless data compression states that the expected codeword length is no less than  $\mathbb{H}(\mathcal{X}^{\chi(\mathbf{u})}) / \log_2 |\mathcal{A}| = w_{\mathbf{H}}(\mathbf{u}) \log_2 |\mathcal{X}| / \log_2 |\mathcal{A}| = L(m, \mathbf{u})$ , and the minimal expected codeword length can be achieved by an optimal source code, e.g., a Huffman code (see [49, Th. 5.4.1]). Here, since the files/sources are i.i.d. according to a uniform distribution, the codeword lengths are identical.

that  $p(\mathbf{q}|m)$  also lies in  $\mathcal{P}_{\mathbf{D}}$ , and hence  $\mathcal{F}_{\mathbf{D}} \subseteq \mathcal{P}_{\mathbf{D}}$ . The intuition behind this fact is that the function  $L(m, \mathbf{u})$  defines the minimum required lengths of answers for a valid single-server WPIR scheme, thus we have more choices of conditional PMFs  $p(\mathbf{q}|m)$  in  $\mathcal{P}_{\mathbf{D}}$ .

Now, if we take the minimization over  $\mathcal{P}_{\mathbf{D}}$ , which is a candidate set that is larger than  $\mathcal{F}_{\mathbf{D}}$ , we have

$$\begin{aligned} \min_{p(\mathbf{q}|m) \in \mathcal{F}_{\mathbf{D}}} I(M; \mathbf{Q}) &\geq \min_{p(\mathbf{q}|m) \in \mathcal{P}_{\mathbf{D}}} I(M; \mathbf{Q}) \\ &\stackrel{(a)}{\geq} \min_{p(\mathbf{u}|m) : \mathbb{E}[L(M, \mathbf{U})] \leq \mathbf{D}} I(M; \mathbf{U}), \end{aligned} \quad (14)$$

where (a) follows directly from the data processing inequality of the first statement of Lemma 1. Therefore, a lower bound to the convex optimization problem (11) is given.

In the second part of the proof, we show that (14) admits a closed-form expression by using a useful result from rate-distortion theory, a lower bound on the rate-distortion function.<sup>6</sup> This lower bound is adapted to (14) and is re-stated as follows.

*Lemma 3* ([47, Th. 9.4.1], [48, Th. 10.19]): Given any  $\mathbf{D} \in [M]$ , we have

$$\begin{aligned} \min_{p(\mathbf{u}|m) : \mathbb{E}[L(M, \mathbf{U})] \leq \mathbf{D}} I(M; \mathbf{U}) \\ \geq \mathbf{H}(M) + \sum_{m \in [M]} P_M(m) \log_2 v_m - \lambda \mathbf{D} \end{aligned} \quad (15)$$

for an arbitrary choice of  $\lambda > 0$  and for any  $v_m, m \in [M]$ , satisfying

$$\sum_{m \in [M]} v_m 2^{-\lambda L(m, \mathbf{u})} \leq 1, \quad \mathbf{u} \in \{0, 1\}^M. \quad (16)$$

We remark again that the length function  $L(m, \mathbf{u})$  is equal to  $w_{\mathbf{H}}(\mathbf{u})$  for all  $m \in \chi(\mathbf{u})$ . By using the condition in (16) for each  $\mathbf{u}$  with  $w_{\mathbf{H}}(\mathbf{u}) = 1$ , we obtain

$$\begin{aligned} \sum_{m \in [M]} v_m 2^{-\lambda L(m, \mathbf{u})} &= \sum_{m \in \chi(\mathbf{u})} v_m 2^{-\lambda w_{\mathbf{H}}(\mathbf{u})} \\ &= v_m 2^{-\lambda} \leq 1, \quad \text{where } m \in \chi(\mathbf{u}). \end{aligned}$$

This implies that for any  $m \in [M]$ ,  $v_m \cdot 2^{-\lambda} \leq 1$ , and hence by symmetry, we can simply assume that

$$v_m = v, \quad \forall m.$$

Next, we apply (16) for all  $\mathbf{u} \in \{0, 1\}^M$ :

$$\begin{aligned} v \cdot 2^{-1\lambda} &\leq 1 && \text{if } w_{\mathbf{H}}(\mathbf{u}) = 1, \\ 2v \cdot 2^{-2\lambda} &\leq 1 && \text{if } w_{\mathbf{H}}(\mathbf{u}) = 2, \\ 3v \cdot 2^{-3\lambda} &\leq 1 && \text{if } w_{\mathbf{H}}(\mathbf{u}) = 3, \\ &&& \vdots \\ Mv \cdot 2^{-M\lambda} &\leq 1 && \text{if } w_{\mathbf{H}}(\mathbf{u}) = M. \end{aligned}$$

<sup>6</sup>The proof of this lower bound is based on the Karush–Kuhn–Tucker optimality conditions, see the details in [47, Ch. 9], [48, Ch. 10].

From the above conditions, we obtain

$$\log_2 v \leq \begin{cases} \lambda & \text{if } \lambda > \log_2 \frac{2}{1}, \\ 2\lambda - \log_2 2 & \text{if } \log_2 \frac{3}{2} < \lambda \leq \log_2 \frac{2}{1}, \\ 3\lambda - \log_2 3 & \text{if } \log_2 \frac{4}{3} < \lambda \leq \log_2 \frac{3}{2}, \\ \vdots & \\ M\lambda - \log_2 M & \text{if } 0 < \lambda \leq \log_2 \frac{M}{M-1}. \end{cases}$$

Now, taking

$$(\log_2 v, \lambda) = \left( w\lambda - \log_2 w, \log_2 \frac{w}{w-1} \right), \quad w \in [2 : M],$$

and substituting this in (15) with  $P_M(m) = \frac{1}{M}$ , we have

$$\begin{aligned} & \min_{p(\mathbf{u}|m): \mathbb{E}[L(M, \mathbf{U})] \leq D} I(M; \mathbf{U}) \\ & \geq \log_2 M + \left( w \log_2 \frac{w}{w-1} - \log_2 w \right) - \left( \log_2 \frac{w}{w-1} \right) D \\ & = \log_2 \frac{M}{w-1} - \left( \log_2 \frac{w}{w-1} \right) (D - (w-1)), \end{aligned} \quad (17)$$

for  $w \in [2 : M]$ .

Here, (17) is a linear function of  $D$  with slope  $-\lambda = -\log_2 \frac{w}{w-1}$ , which is strictly increasing in  $w \in [2 : M]$ . Therefore, the best lower bound for (14) is the piecewise function

$$\rho_{\text{LB}}^{(M)}(D) = \begin{cases} \log_2 \frac{M}{1} - \left( \log_2 \frac{2}{1} \right) (D - 1) & \text{if } 1 \leq D \leq 2, \\ \log_2 \frac{M}{2} - \left( \log_2 \frac{3}{2} \right) (D - 2) & \text{if } 2 < D \leq 3, \\ \vdots & \\ \log_2 \frac{M}{M-1} - \left( \log_2 \frac{M}{M-1} \right) (D - (M-1)) & \text{if } M-1 < D \leq M. \end{cases}$$

See also the pictorial illustration in Fig. 3. For instance, the red line  $\log_2 \frac{M}{w-1} - \left( \log_2 \frac{w}{w-1} \right) (D - (w-1))$  going through the points  $(w-1, \log_2 \frac{M}{w-1})$  and  $(w, \log_2 \frac{M}{w})$  has the largest function values over the interval  $[w-1, w]$ .

Since the leakage-download function of an arbitrary WPIR scheme is bounded from below by  $\rho_{\text{LB}}^{(M)}(D)$ , and it can be shown that the pair  $(\bar{\rho}, D^{(M)}(\bar{\rho}))$  of Theorem 2 lies on  $\rho_{\text{LB}}^{(M)}(D)$  (details are omitted for brevity), this completes the converse proof.

## VII. CONVERSE OF THEOREM 3

Following an argumentation similar to the first part of the converse proof of Theorem 2, since the MaxL metric also satisfies the data processing inequality from the first item of Lemma 1, the leakage-download function for the MaxL privacy metric can be bounded from below by

$$\min_{p(\mathbf{u}|m): \mathbb{E}[L(M, \mathbf{U})] \leq D} \text{MaxL}(M; \mathbf{U}). \quad (18)$$

This by itself is not a convex minimization problem. In this proof, we derive a lower bound to (18) directly. To make the problem tractable, we use the fact that  $2^{\text{MaxL}(M; \mathbf{U})}$  is convex in  $P_{\mathbf{U}|M}$  (see 2) in Lemma 1).

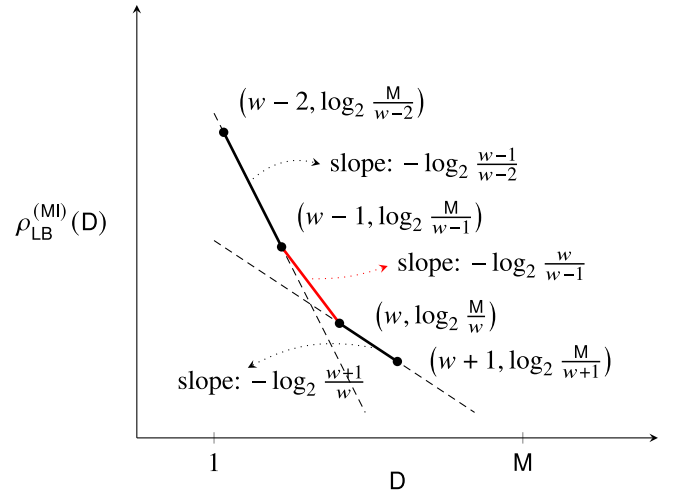


Fig. 3. Illustration of the function  $\rho_{\text{LB}}^{(M)}(D)$ , which is defined by many linear functions. The slope of the linear functions is strictly increasing in  $w \in [2 : M]$ .

We know from (3) that maximizing the objective function  $\text{MaxL}(M; \mathbf{U})$  is equivalent to maximizing the function

$$2^{\text{MaxL}(M; \mathbf{U})} = \sum_{\mathbf{u} \in \{0,1\}^M} \max_{m \in [M]} p(\mathbf{u}|m).$$

Moreover, from (13) we know that

$$p(\mathbf{u}|m) = 0, \quad \text{if } \mathbf{u} \in \{0,1\}^M, m \notin \chi(\mathbf{u}).$$

Thus, (18) can be re-written as the convex minimization problem

$$\text{minimize } \sum_{\mathbf{u} \in \{0,1\}^M} \max_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \quad (19a)$$

$$\text{subject to } \sum_{\mathbf{u} \in \{0,1\}^M, m \in [M]} P_M(m) p(\mathbf{u}|m) L(m, \mathbf{u}) \leq D, \quad (19b)$$

$$\sum_{\mathbf{u} \in \{0,1\}^M} p(\mathbf{u}|m) = 1, \quad \forall m \in [M]. \quad (19c)$$

Furthermore, using the fact that

$$\sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \leq w_H(\mathbf{u}) \max_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m), \quad \forall \mathbf{u} \in \{0,1\}^M,$$

and

$$p(\mathbf{1}_{\{m\}}|m) = 1 - \sum_{\substack{\mathbf{u}: m \in \chi(\mathbf{u}) \\ w_H(\mathbf{u}) > 1}} p(\mathbf{u}|m), \quad \forall m \in [M], \quad (20)$$

the objective function (19a) becomes<sup>7</sup>

$$\begin{aligned} & \sum_{\mathbf{u}} \max_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \\ & = \sum_{\mathbf{u}: w_H(\mathbf{u})=1} p(\mathbf{u}|\chi(\mathbf{u})) + \sum_{\mathbf{u}: w_H(\mathbf{u})>1} \max_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \end{aligned}$$

<sup>7</sup>In the following, the ranges of the summations and also the explicit summation variable are sometimes omitted as they are clear from the context.

TABLE II  
EQUIVALENT LP PROBLEM OF (21) WITH THE OBJECTIVE FUNCTION  $\mathbf{c}[\mathbf{y}, \mathbf{z}]^T = -\sum_{w=2}^M (1 - \frac{1}{w})y_w$ , WHERE THE FIRST ROW INDICATES THE VARIABLES  $[\mathbf{y}, \mathbf{z}] = [y_2, \dots, y_M, z_1, z_2]$ , THE SECOND ROW REPRESENTS THE COEFFICIENTS  $\mathbf{c}$ , AND THE THIRD AND FOURTH ROWS ARE OBTAINED FROM THE CONSTRAINTS  $\mathbf{A}[\mathbf{y}, \mathbf{z}]^T = \mathbf{b}^T$

	$y_2$	$y_3$	$\dots$	$y_{w-1}$	$y_w$	$\dots$	$y_M$	$z_1$	$z_2$	$\mathbf{b}$
$\mathbf{c}$	$-1 - \frac{1}{2}$	$-1 - \frac{1}{3}$	$\dots$	$-1 - \frac{1}{w-1}$	$-1 - \frac{1}{w}$	$\dots$	$-1 - \frac{1}{M}$	0	0	
$\mathbf{A}$	1	1	$\dots$	1	1	$\dots$	1	1	0	M
	1	2	$\dots$	$w-2$	$w-1$	$\dots$	$M-1$	0	1	$M(D-1)$

$$\begin{aligned}
 &= \sum_{m \in [M]} \left[ 1 - \sum_{\substack{\mathbf{u}: m \in \chi(\mathbf{u}) \\ w_H(\mathbf{u}) > 1}} p(\mathbf{u}|m) \right] + \sum_{\mathbf{u}: w_H(\mathbf{u}) > 1} \max_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \\
 &\geq M - \sum_{m \in [M]} \sum_{\substack{\mathbf{u}: m \in \chi(\mathbf{u}) \\ w_H(\mathbf{u}) > 1}} p(\mathbf{u}|m) + \sum_{\mathbf{u}: w_H(\mathbf{u}) > 1} \sum_{m \in \chi(\mathbf{u})} \frac{p(\mathbf{u}|m)}{w_H(\mathbf{u})} \\
 &\stackrel{(a)}{=} M - \sum_{\mathbf{u}: w_H(\mathbf{u}) > 1} \sum_{m \in \chi(\mathbf{u})} \left( 1 - \frac{1}{w_H(\mathbf{u})} \right) p(\mathbf{u}|m),
 \end{aligned}$$

where (a) holds since by expanding the double summation, one can see that

$$\sum_{m \in [M]} \sum_{\substack{\mathbf{u}: m \in \chi(\mathbf{u}) \\ w_H(\mathbf{u}) > 1}} p(\mathbf{u}|m) = \sum_{\mathbf{u}: w_H(\mathbf{u}) > 1} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m).$$

Similarly, by substituting (20) into the download cost constraint (19c), we get

$$\begin{aligned}
 &\sum_{\mathbf{u}} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) L(m, \mathbf{u}) \\
 &= \sum_{w_H(\mathbf{u})=1} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \cdot 1 \\
 &\quad + \sum_{w_H(\mathbf{u}) > 1} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) w_H(\mathbf{u}) \\
 &= M - \sum_{w_H(\mathbf{u}) > 1} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \\
 &\quad + \sum_{w_H(\mathbf{u}) > 1} w_H(\mathbf{u}) \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \\
 &= M + \sum_{w_H(\mathbf{u}) > 1} (w_H(\mathbf{u}) - 1) \left[ \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \right] \leq MD.
 \end{aligned}$$

Next, define  $y_w \triangleq \sum_{\mathbf{u}: w_H(\mathbf{u})=w} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m)$  for  $w \in [2 : M]$ . Because

$$\begin{aligned}
 &\sum_{\mathbf{u}: w_H(\mathbf{u}) > 1} (w_H(\mathbf{u}) - 1) \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m) \\
 &= \sum_{w=2}^M (w-1) \sum_{\mathbf{u}: w_H(\mathbf{u})=w} \sum_{m \in \chi(\mathbf{u})} p(\mathbf{u}|m),
 \end{aligned}$$

it can be shown that a lower bound to (19) can be computed from the linear programming (LP) formulation

$$\text{minimize } \rho_{\text{LB}}^{(\text{MaxL})}(\mathcal{D}) \triangleq M - \sum_{w=2}^M \left[ 1 - \frac{1}{w} \right] y_w \quad (21a)$$

$$\text{subject to } \sum_{w=2}^M y_w \leq M, \quad (21b)$$

$$\sum_{w=2}^M (w-1) \cdot y_w \leq M(D-1), \quad (21c)$$

with variables  $y_w$ ,  $w \in [2 : M]$ .

We convert the inequalities in the constraints (21b) and (21c) to equalities by introducing variables  $z_1$  and  $z_2$ . Thus, we have the constraints as

$$\sum_{w=2}^M y_w + z_1 = M,$$

$$\sum_{w=2}^M (w-1) \cdot y_w + z_2 = M(D-1).$$

Now, define  $c_w \triangleq -(1 - \frac{1}{w})$ , for  $w \in [2 : M]$ , and let  $c_{M+1} = c_{M+2} = 0$ . Then, the objective function of (21) can be written in matrix form as  $\mathbf{M} + \mathbf{c}[y_2, \dots, y_M, z_1, z_2]^T = M - \sum_{w=2}^M (1 - \frac{1}{w})y_w$ , with  $\mathbf{c} = (c_2, \dots, c_{M+2})$ , and the constraints as  $\mathbf{A}[\mathbf{y}, \mathbf{z}]^T = \mathbf{b}^T$ , where

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & \dots & 1 & 1 & 0 \\ 1 & 2 & \dots & M-1 & 0 & 1 \end{pmatrix}, \quad \mathbf{b}^T = \begin{pmatrix} M \\ M(D-1) \end{pmatrix},$$

$\mathbf{y} = (y_2, \dots, y_M)$ , and  $\mathbf{z} = (z_1, z_2)$ . The equivalent LP problem of (21) (without the constant  $M$ ) is shown in Table II.

Consider the standard LP problem of minimizing  $\mathbf{c}\mathbf{x}^T$  subject to  $\mathbf{A}\mathbf{x}^T = \mathbf{b}^T$  and  $\mathbf{x} \geq \mathbf{0}$ , where  $\mathbf{A}$  is an arbitrary  $m \times n$  matrix. A *basic solution* is any solution where a subset of  $n-m$  variables are zero. The  $m$  nonzero variables of a basic solution are referred to as the basic variables, while the remaining variables are known as the nonbasic variables. Let  $\mathbf{A}_{\mathcal{B}}$  denote the submatrix of  $\mathbf{A}$  consisting of the  $m$  columns of  $\mathbf{A}$  corresponding to the basic variables, and by  $\mathbf{A}_{\mathcal{N}}$  the submatrix consisting of the remaining  $n-m$  columns corresponding to the nonbasic variables. Denote by  $\mathbf{x}_{\mathcal{B}}$  the basic variables of a basic solution  $\mathbf{x}$ , and let  $\mathbf{c}_{\mathcal{B}}$  and  $\mathbf{c}_{\mathcal{N}}$  be the subvectors of  $\mathbf{c}$  that correspond to the basic and nonbasic variables, respectively.

The following proposition provides a sufficient condition for a basic solution to be optimal for an LP problem.<sup>8</sup>

*Proposition 1 ([50, p. 44]):* If there exists a basic solution such that the *relative cost vector* for nonbasic variables, defined as  $\mathbf{c}_{\mathcal{N}} - \mathbf{c}_{\mathcal{B}}\mathbf{A}_{\mathcal{B}}^{-1}\mathbf{A}_{\mathcal{N}}$ , is nonnegative, then the basic solution is optimal.

From the capacity-achieving scheme proposed in Section V, one can show that

$$p^*(\mathbf{u}|m) = \begin{cases} \frac{\frac{w-D}{(M-1)}}{\binom{w-2}{M-1}} & \text{if } w_{\mathcal{H}}(\mathbf{u}) = w-1, \\ \frac{D-(w-1)}{\binom{M-1}{w-1}} & \text{if } w_{\mathcal{H}}(\mathbf{u}) = w, \\ 0 & \text{otherwise,} \end{cases}$$

for  $w-1 \leq D \leq w$ ,  $w \in [2 : M]$ , where the superscript  $*$  indicates that the corresponding quantity is for the particular scheme from Section V. Then, since  $y_w^* = \sum_{\mathbf{u}:w_{\mathcal{H}}(\mathbf{u})=w} \sum_{m \in \chi(\mathbf{u})} p^*(\mathbf{u}|m)$  for  $w \in [2 : M]$ , we have, for  $1 \leq D \leq 2$ ,

$$\begin{aligned} y_2^* &= \sum_{\mathbf{u}:w_{\mathcal{H}}(\mathbf{u})=2} \sum_{m \in \chi(\mathbf{u})} \frac{D-1}{M-1} = M(D-1), \\ y_{w'}^* &= 0, \quad w' \in [3 : M], \\ z_1^* &= M(2-D), \quad z_2^* = 0, \end{aligned}$$

and

$$\begin{aligned} y_{w-1}^* &= \binom{M}{w-1} (w-1) \frac{w-D}{\binom{M-1}{w-2}} = M(w-D), \\ y_w^* &= \binom{M}{w} w \frac{D-(w-1)}{\binom{M-1}{w-1}} = M(D-(w-1)), \\ y_{w'}^* &= 0, \quad w' \in [3 : M] \setminus \{w-1, w\}, \\ z_1^* &= 0, \quad z_2^* = 0, \end{aligned}$$

for  $w-1 \leq D \leq w$ ,  $w \in [3 : M]$ . To prove the converse part, we use Proposition 1 to verify the optimality of  $y_w^*$ ,  $w \in [2 : M]$ , and  $z_1^*$ ,  $z_2^*$ , for (21), by considering two special cases as follows.

*Case 1 ( $1 \leq D \leq 2$ ):* In this case,  $y_2$  and  $z_1$  are the basic variables. Since

$$\begin{aligned} c_w - \mathbf{c}_{\mathcal{B}}\mathbf{A}_{\mathcal{B}}^{-1}\mathbf{A}_{\{y_w\}} &= -\left[1 - \frac{1}{w}\right] - \left(-\frac{1}{2}, 0\right) \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ w-1 \end{pmatrix} \\ &= \frac{w-1}{2} - \frac{w-1}{w} > 0, \end{aligned}$$

for  $w \in [3 : M]$ , and  $0 - \left(-\frac{1}{2}, 0\right) \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}^{-1} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \frac{1}{2} > 0$  for  $z_2$ , it follows that the basic solution is optimal for  $M \geq 2$ .

*Case 2 ( $w-1 < D \leq w$ ,  $w \in [3 : M]$ ):* In this case,  $y_{w-1}$  and  $y_w$  are the basic variables. We can see that for  $w' \in [3 : M] \setminus \{w-1, w\}$ , the relative cost coefficient for  $y_{w'}$  is

$$\begin{aligned} c_{w'} - \mathbf{c}_{\mathcal{B}}\mathbf{A}_{\mathcal{B}}^{-1}\mathbf{A}_{\{y_{w'}\}} &= -\left[1 - \frac{1}{w'}\right] - \left(-\left[1 - \frac{1}{w-1}\right], -\left[1 - \frac{1}{w}\right]\right) \\ &\quad \times \begin{pmatrix} 1 & 1 \\ w-2 & w-1 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ w'-1 \end{pmatrix} \end{aligned}$$

<sup>8</sup>The proof of Proposition 1 can be found in the LP literature, see [50, pp. 43–44].

$$= \frac{1}{w'} - \frac{w-w'}{w-1} + \frac{w-w'-1}{w}.$$

By performing some simple calculations, we get

$$\begin{aligned} \frac{1}{w'} - \frac{w-w'}{w-1} + \frac{w-w'-1}{w} \\ = (w-w') \left[ \frac{w-(w'+1)}{w \cdot w'(w-1)} \right] > 0, \end{aligned}$$

either for  $w' > w-1$  or  $w' < w$ . Moreover, it is easy to see that the relative cost coefficients for  $z_1$  and  $z_2$  are  $\frac{w-2}{w} > 0$  and  $\frac{1}{w-1} - \frac{1}{w} > 0$ , respectively, which implies that the given basic solution is optimal.

## VIII. CONCLUSION

We considered relaxing the perfect privacy condition in single-server PIR and presented a scheme for the studied weakly-private scenario, referred to as WPIR. In doing so, we showed that one can trade off privacy to gain in terms of download cost. Furthermore, we characterized the information leaked using two different metrics: MI and MaxL. The latter is known to be a more robust metric to measure information leakage. Finally, we derived the single-server WPIR capacity for both the MI and MaxL metrics, and showed that the proposed protocol is capacity-achieving. As a final note, we drew the connection between WPIR and rate-distortion theory.

An interesting direction for future work is the derivation of fundamental bounds on other performance metrics like the upload cost and the access complexity for the single server scenario.

## ACKNOWLEDGMENT

The authors would like to thank the three anonymous reviewers and the Guest Editor Prof. Lalitha Sankar for their valuable and insightful comments.

## REFERENCES

- [1] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan, "Private information retrieval," in *Proc. 36th Annu. IEEE Symp. Found. Comput. Sci. (FOCS)*, Milwaukee, WI, USA, Oct. 1995, pp. 41–50.
- [2] N. B. Shah, K. V. Rashmi, and K. Ramchandran, "One extra bit of download ensures perfectly private information retrieval," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Honolulu, HI, USA, Jun./Jul. 2014, pp. 856–860.
- [3] T. H. Chan, S.-W. Ho, and H. Yamamoto, "Private information retrieval for coded storage," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Hong Kong, Jun. 2015, pp. 2842–2846.
- [4] 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.
- [5] R. Tajeddine, O. W. Gnilke, and S. El Rouayheb, "Private information retrieval from MDS coded data in distributed storage systems," *IEEE Trans. Inf. Theory*, vol. 64, no. 11, pp. 7081–7093, Nov. 2018.
- [6] K. Banawan and S. Ulukus, "The capacity of private information retrieval from coded databases," *IEEE Trans. Inf. Theory*, vol. 64, no. 3, pp. 1945–1956, Mar. 2018.
- [7] S. Kumar, H.-Y. Lin, E. Rosnes, and A. Graell i Amat, "Achieving maximum distance separable private information retrieval capacity with linear codes," *IEEE Trans. Inf. Theory*, vol. 65, no. 7, pp. 4243–4273, Jul. 2019.
- [8] H.-Y. Lin, S. Kumar, E. Rosnes, and A. Graell i Amat, "Asymmetry helps: Improved private information retrieval protocols for distributed storage," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Guangzhou, China, Nov. 2018, pp. 1–5.

- [9] R. Freij-Hollanti, O. W. Gnilke, C. Hollanti, A.-L. Horlemann-Trautmann, D. Karpuk, and I. Kubjas, “ $t$ -private information retrieval schemes using transitive codes,” *IEEE Trans. Inf. Theory*, vol. 65, no. 4, pp. 2107–2118, Apr. 2019.
- [10] H. Sun and S. A. Jafar, “The capacity of robust private information retrieval with colluding databases,” *IEEE Trans. Inf. Theory*, vol. 64, no. 4, pp. 2361–2370, Apr. 2018.
- [11] R. Freij-Hollanti, O. W. Gnilke, C. Hollanti, and D. A. Karpuk, “Private information retrieval from coded databases with colluding servers,” *SIAM J. Appl. Algebra Geom.*, vol. 1, no. 1, pp. 647–664, Nov. 2017.
- [12] R. G. L. D’Oliveira and S. El Rouayheb, “One-shot PIR: Refinement and lifting,” *IEEE Trans. Inf. Theory*, vol. 66, no. 4, pp. 2443–2455, Apr. 2020.
- [13] L. Holzbaur, R. Freij-Hollanti, J. Li, and C. Hollanti. (Mar. 2019). *Towards the Capacity of Private Information Retrieval From Coded and Colluding Servers*. [Online]. Available: <https://arxiv.org/abs/1903.12552>
- [14] K. Banawan and S. Ulukus, “The capacity of private information retrieval from Byzantine and colluding databases,” *IEEE Trans. Inf. Theory*, vol. 65, no. 2, pp. 1206–1219, Feb. 2019.
- [15] C. Tian, H. Sun, and J. Chen, “Capacity-achieving private information retrieval codes with optimal message size and upload cost,” *IEEE Trans. Inf. Theory*, vol. 65, no. 11, pp. 7613–7627, Nov. 2019.
- [16] Y. Zhang, E. Yaakobi, T. Etzion, and M. Schwartz, “On the access complexity of PIR schemes,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Paris, France, Jul. 2019, pp. 2134–2138.
- [17] S. Kadhe, B. Garcia, A. Heidarzadeh, S. El Rouayheb, and A. Sprintson, “Private information retrieval with side information,” *IEEE Trans. Inf. Theory*, vol. 66, no. 4, pp. 2032–2043, Apr. 2020.
- [18] A. Heidarzadeh, F. Kazemi, and A. Sprintson, “The role of coded side information in single-server private information retrieval,” *IEEE Trans. Inf. Theory*, vol. 67, no. 1, pp. 25–44, Jan. 2021.
- [19] N. Raviv, I. Tamo, and E. Yaakobi, “Private information retrieval in graph-based replication systems,” *IEEE Trans. Inf. Theory*, vol. 66, no. 6, pp. 3590–3602, Jun. 2020.
- [20] H. Yang, W. Shin, and J. Lee, “Private information retrieval for secure distributed storage systems,” *IEEE Trans. Inf. Forensics Security*, vol. 13, no. 12, pp. 2953–2964, Dec. 2018.
- [21] Z. Jia, H. Sun, and S. A. Jafar, “Cross subspace alignment and the asymptotic capacity of  $X$ -secure  $T$ -private information retrieval,” *IEEE Trans. Inf. Theory*, vol. 65, no. 9, pp. 5783–5798, Sep. 2019.
- [22] Y. Zhang, T. Etzion, and E. Yaakobi, “Bounds on the length of functional PIR and batch codes,” *IEEE Trans. Inf. Theory*, vol. 66, no. 8, pp. 4917–4934, Aug. 2020.
- [23] Y. Zhang, E. Yaakobi, and T. Etzion. (Jul. 2019). *Private Proximity Retrieval Codes*. [Online]. Available: <https://arxiv.org/abs/1907.10724>
- [24] H.-Y. Lin, S. Kumar, E. Rosnes, A. Graell i Amat, and E. Yaakobi, “Weakly-private information retrieval,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Paris, France, Jul. 2019, pp. 1257–1261.
- [25] I. Samy, R. Tandon, and L. Lazos, “On the capacity of leaky private information retrieval,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Paris, France, Jul. 2019, pp. 1262–1266.
- [26] H.-Y. Lin, S. Kumar, E. Rosnes, A. Graell i Amat, and E. Yaakobi. (Jul. 2020). *Multi-Server Weakly-Private Information Retrieval*. [Online]. Available: <https://arxiv.org/abs/2007.10174>
- [27] K. Bhattad and K. R. Narayanan, “Weakly secure network coding,” in *Proc. Int. Symp. Netw. Coding (NetCod)*, Riva del Garda, Italy, Apr. 2005.
- [28] D. Silva and F. R. Kschischang, “Universal weakly secure network coding,” in *Proc. IEEE Inf. Theory Workshop (ITW)*, Volos, Greece, Jun. 2009, pp. 281–285.
- [29] M. Bloch and J. Barros, *Physical-Layer Security: From Information Theory to Security Engineering*. Cambridge, U.K.: Cambridge Univ. Press, 2011.
- [30] B. Köpf and D. Basin, “An information-theoretic model for adaptive side-channel attacks,” in *Proc. 14th ACM Conf. Comput. Commun. Security (CCS)*, Alexandria, VA, USA, Oct./Nov. 2007, pp. 286–296.
- [31] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Proc. Theory Cryptography Conf. (TCC)*, New York, NY, USA, Mar. 2006, pp. 265–284.
- [32] C. Dwork, “Differential privacy,” in *Proc. 33rd Int. Colloquium Automata Lang. Program. (ICALP)*, Venice, Italy, Jul. 2006, pp. 1–12.
- [33] R. Zhou, T. Guo, and C. Tian, “Weakly private information retrieval under the maximal leakage metric,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Los Angeles, CA, USA, Jun. 2020, pp. 1089–1094.
- [34] D. Asonov and J.-C. Freytag, “Repudiative information retrieval,” in *Proc. ACM Workshop Privacy Electron. Soc. (WPES)*, Washington, DC, USA, Nov. 2002, pp. 32–40.
- [35] R. R. Toledo, G. Danezis, and I. Goldberg, “Lower-cost  $\epsilon$ -private information retrieval,” in *Proc. Privacy Enhancing Technol. Symp. (PETS)*, Darmstadt, Germany, Jul. 2016, pp. 184–201.
- [36] T. Guo, R. Zhou, and C. Tian, “On the information leakage in private information retrieval systems,” *IEEE Trans. Inf. Forensics Security*, vol. 15, pp. 2999–3012, Mar. 2020.
- [37] I. Samy, M. A. Attia, R. Tandon, and L. Lazos. (Jun. 2020). *Asymmetric Leaky Private Information Retrieval*. [Online]. Available: <https://arxiv.org/abs/2006.03048>
- [38] H. Sun and S. A. Jafar, “The capacity of symmetric private information retrieval,” *IEEE Trans. Inf. Theory*, vol. 65, no. 1, pp. 322–329, Jan. 2019.
- [39] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, “Private information retrieval,” *J. ACM*, vol. 45, no. 6, pp. 965–982, Nov. 1998.
- [40] I. Samy, M. A. Attia, R. Tandon, and L. Lazos, “Latent-variable private information retrieval,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Los Angeles, CA, USA, Jun. 2020, pp. 1071–1076.
- [41] G. Smith, “On the foundations of quantitative information flow,” in *Proc. 12th Int. Conf. Found. Softw. Sci. Comput. Struct. (FoSSaCS)*, York, U.K., Mar. 2009, pp. 288–302.
- [42] G. Barthe and B. Köpf, “Information-theoretic bounds for differentially private mechanisms,” in *Proc. 24th IEEE Comput. Security Found. Symp. (CSF)*, Cernay-la-Ville, France, Jun. 2011, pp. 191–204.
- [43] I. Issa, A. B. Wagner, and S. Kamath, “An operational approach to information leakage,” *IEEE Trans. Inf. Theory*, vol. 66, no. 3, pp. 1625–1657, Mar. 2020.
- [44] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, “What can we learn privately?” in *Proc. 49th Annu. IEEE Symp. Found. Comput. Sci. (FOCS)*, Philadelphia, PA, USA, Oct. 2008, pp. 531–540.
- [45] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, “Local privacy and statistical minimax rates,” in *Proc. 54th Annu. IEEE Symp. Found. Comput. Sci. (FOCS)*, Berkeley, CA, USA, Oct. 2013, pp. 429–438.
- [46] S. Vithana, K. Banawan, and S. Ulukus. (Mar. 2020). *Semantic Private Information Retrieval*. [Online]. Available: <https://arxiv.org/abs/2003.13667>
- [47] R. G. Gallager, *Information Theory and Reliable Communication*. New York, NY, USA: Wiley, 1968.
- [48] S. M. Moser, *Advanced Topics in Information Theory (Lecture Notes)*, Signal Inf. Process. Lab., ETH Zürich, Switzerland, and Dept. Elect. Comput. Eng., Nat. Chiao Tung Univ., Hsinchu, Taiwan, 2019. [Online]. Available: <http://moser-isi.ethz.ch/scripts.html>
- [49] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New York, NY, USA: Wiley, 2006.
- [50] D. G. Luenberger and Y. Ye, *Linear and Nonlinear Programming*, 4th ed. Cham, Switzerland: Springer Int., 2016.