

The Capacity of Single-Server Weakly-Private Information Retrieval

Hsuan-Yin Lin*, Siddhartha Kumar*, Eirik Rosnes*, Alexandre Graell i Amat^{†*}, and Eitan Yaakobi[‡]

*Simula UiB, N-5006 Bergen, Norway

[†]Department of Electrical Engineering, Chalmers University of Technology, SE-41296 Gothenburg, Sweden

[‡]Department of Computer Science, Technion — Israel Institute of Technology, Haifa, 3200003 Israel

Abstract—Weakly-private information retrieval (WPIR) is a variant of the private information retrieval problem in which a user wants to efficiently retrieve a file stored across a set of servers while tolerating some information leakage on the identity of the requested file to the servers. In this paper, we consider WPIR from a single-server database where the information leakage is measured in terms of the mutual information (MI) or maximal leakage (MaxL) privacy metrics. In particular, we establish a connection between the WPIR problem and rate-distortion theory, and fully characterize the optimal tradeoff between the download cost and the allowed information leakage under the MI and MaxL metrics, settling the single-server WPIR capacity.

I. INTRODUCTION

With the ever increasing demand for privacy rights on the Internet, preserving user privacy has become of vital importance for, e.g., big data applications and distributed information systems. The European Union’s General Data Protection Regulation is an important example that demonstrates the increasing awareness of this issue.

Private information retrieval (PIR), introduced by Chor *et al.* [1], is one of the fundamental primitives that ensures user privacy. In particular, it enables a user to download a desired file stored in one or several servers without disclosing the identity of the requested file to any of the servers. In recent years, designing efficient PIR schemes from an information-theoretic perspective has attracted significant attention [2], [3]. In this line of research, the efficiency is mainly measured in terms of download rate, defined as the ratio between the requested file size and the expected number of downloaded symbols. The maximum possible download rate is referred to as the *PIR capacity*.

Characterizing the capacity for variants of the PIR problem is of fundamental importance. Sun and Jafar derived the PIR capacity for the classical case where data is replicated across multiple servers [2]. For the case when data is encoded by a maximum distance separable (MDS) code and then stored in multiple servers, the capacity, referred to as the *MDS-PIR capacity*, was given in [3]. The PIR capacity when the storage code is from a particular family of non-MDS codes, introduced

in [4], was determined and shown to be equal to the MDS-PIR capacity in [5]. Other PIR capacity results for various scenarios can be found in, e.g., [6]–[9].

Recently, the notion of weakly-private information retrieval (WPIR) was introduced independently by the works [10]–[12]. WPIR can be seen as a generalization of PIR, which guarantees the retrieval of a single file while leaking partial information about the identity of the requested file to the servers. For the scenario with multiple noncolluding servers, it was shown that the download rate as well as the upload cost and access complexity can be improved by trading off the information leakage.

In practice, the servers are typically in control of a single provider, violating the assumption that they can not collude. Motivated by this important and practical observation, in this paper we turn our attention and initiate the study of the more compelling scenario in which data is stored in a single server, i.e., the problem of single-server WPIR is considered. Similar to [11], to measure the information leakage in an information-theoretic sense, we adopt two commonly-used information leakage metrics in the information theory literature, the mutual information (MI) and maximal leakage (MaxL) metrics [13]–[15]. The latter is known as one of the most useful information-theoretic measures of information leakage in the computer security literature. In particular, by developing the similarities between the WPIR problem and rate-distortion theory, we fully characterize the optimal tradeoff between the download cost and a given level of information leakage for the MI and MaxL privacy metrics, and hence determine the single-server WPIR capacity. Moreover, a novel single-server WPIR capacity-achieving scheme for both metrics is constructed by using a time-sharing approach. Due to lack of space, all technical proofs are omitted, and we refer to the extended version of this work [16] for full details as well as the converse proofs.

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 and sets by calligraphic uppercase letters, e.g., \mathbf{x} and \mathcal{X} , respectively. In general, vectors are represented as row vectors throughout the paper. We use uppercase letters for random variables (RVs)

This work was partially funded by the Swedish Research Council (grant 2016-04253), the Israel Science Foundation (grant #1817/18), and by the Technion Hiroshi Fujiwara Cyber Security Research Center and the Israel National Cyber Directorate.

(either scalar or vector), e.g., X or \mathbf{X} . For a given index set \mathcal{S} , we write $X^{\mathcal{S}}$ to represent $\{X^{(m)} : m \in \mathcal{S}\}$. 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. $H(X)$, $H(P_X)$, or $H(p_1, \dots, p_{|\mathcal{X}|})$ 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 . The Galois field with q elements is denoted by $\text{GF}(q)$.

B. System Model

We consider a single server that stores M independent files $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$, where each file $\mathbf{X}^{(m)} = (X_1^{(m)}, \dots, X_\beta^{(m)})$, $m \in [M]$, is represented as a length- β row vector over $\text{GF}(q)$. Assume that each element of $\mathbf{X}^{(m)}$ is chosen independently and uniformly at random from $\text{GF}(q)$. Thus, in q -ary units, we have $H(\mathbf{X}^{(m)}) = \beta$, $\forall m \in [M]$. A user wishes to efficiently retrieve $\mathbf{X}^{(M)}$ by allowing some information leakage to the server, where the requested file index M is assumed to be uniformly distributed over $[M]$. 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 \mathcal{S} , whose alphabet is \mathcal{S} .
- A query function

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

that generates a query $\mathbf{Q} = \phi(M, \mathcal{S})$ with alphabet \mathcal{Q} . The query \mathbf{Q} is sent to the server to retrieve the M -th file.

- An answer function

$$\varphi: \mathcal{Q} \times \text{GF}(q)^{\beta M} \rightarrow \mathcal{A}^{\beta L}$$

that returns the answer $\mathbf{A} \triangleq \varphi(\mathbf{Q}, \mathbf{X}^{[M]})$ 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} .

- A privacy leakage metric $\rho^{(\cdot)}(P_{\mathbf{Q}|M}) \geq 0$, which is defined as a function of the conditional probability mass function (PMF) $P_{\mathbf{Q}|M}$, 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.

This scheme must satisfy the condition of perfect retrievability,

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

We remark that a PIR scheme is equivalent 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 , its designed query conditional PMF 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 paper, 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, can be defined based on the min-entropy (MinE) measure discussed in the computer science literature, see, e.g., [13]. Moreover, since we assume that M is uniformly distributed, the MinE information leakage and the MaxL privacy metric can be shown to be equivalent [14], [15]. The latter is given by

$$\text{MaxL}(M; \mathbf{Q}) \triangleq \log_2 \sum_{\mathbf{q} \in \mathcal{Q}} \max_{m \in [M]} P_{\mathbf{Q}|M}(\mathbf{q}|m). \quad (3)$$

Henceforth, we will use the closed-form expression in (3) as the MaxL privacy metric and denote the MaxL privacy metric of a WPIR scheme by

$$\rho^{(\text{MaxL})}(P_{\mathbf{Q}|M}) \triangleq \text{MaxL}(M; \mathbf{Q}).$$

In the following, the information leakage metric of a WPIR scheme \mathcal{C} is denoted by $\rho^{(\cdot)}(\mathcal{C})$.

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

For WPIR, in contrast to PIR, the download costs for the retrieval of different files do not necessarily need to be the same. 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 expected length of the returned answer over all random queries,

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

and the overall download cost is measured in terms of the expected download cost over all files, i.e.,

$$D(\mathcal{C}) \triangleq \mathbb{E}_{P_M}[\mathbb{E}_{P_{\mathbf{Q}|M}}[L(\mathbf{Q})]] = \mathbb{E}_{M, \mathbf{Q}}[L(\mathbf{Q})]. \quad (4)$$

Accordingly, the WPIR rate is defined as $R(\mathcal{C}) \triangleq D(\mathcal{C})^{-1}$.

In this paper, our goal is to characterize the optimal trade-off 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, ϱ) is said to be achievable in terms of the information leakage metric $\rho^{(\cdot)}$ if there exists a WPIR scheme \mathcal{C} such that $\mathbb{E}_{M, \mathbf{Q}}[L(\mathbf{Q})] \leq D$ and $\rho^{(\cdot)}(\mathcal{C}) \leq \varrho$. The download-leakage region is the set of all achievable download-leakage pairs (D, ϱ) .

III. CHARACTERIZATION OF THE OPTIMAL DOWNLOAD-LEAKAGE TRADEOFF

Consider a single-server WPIR scheme, where the leakage is measured by $\rho^{(\text{MI})}$ or $\rho^{(\text{MaxL})}$. The minimum achievable download cost for a given leakage constraint ϱ can be formulated by the optimization problem

$$\begin{aligned} & \text{minimize} && \mathbb{E}_{P_M P_{Q|M}}[L(\mathbf{Q})] \\ & \text{subject to} && P_{Q|M} \in \mathcal{P}_{\text{ret}}, \end{aligned} \quad (5a)$$

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

where \mathcal{P}_{ret} is defined as the set of all PMFs that satisfy (1).

For convenience, since we know that a designed conditional distribution $P_{Q|M}$ of a WPIR scheme always satisfies (5a), throughout this paper we will assume that any $P_{Q|M} \in \mathcal{P}_{\text{ret}}$.

A. The 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. *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 ϱ such that (D, ϱ) is achievable, i.e.,*

$$D^{(\cdot)}(\varrho) \triangleq \min_{P_{Q|M}: \rho^{(\cdot)}(P_{Q|M}) \leq \varrho} \mathbb{E}_{P_M P_{Q|M}}[L(\mathbf{Q})].$$

Accordingly, the single-server WPIR capacity is $C^{(\cdot)}(\varrho) = [D^{(\cdot)}(\varrho)]^{-1}$.

It is known that the single-server PIR capacity is $C_M = \frac{1}{M}$ [1].

Definition 4. *The leakage-download function $\rho^{(\cdot)}(D)$ for single-server WPIR is the minimum of all possible information leakages ϱ for a given download cost constraint D such that (D, ϱ) is achievable.*

Lemma 1. *The MI download-leakage function*

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

is convex in ϱ , while the MaxL download-leakage function

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

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

The convexity of $D^{(\text{MI})}$ 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 (D_1, ϱ_1) and (D_2, ϱ_2) are achievable. Then, for any $\lambda \in [0, 1]$, the pair $(D_\lambda = (1-\lambda)D_1 + \lambda D_2, \varrho_\lambda = (1-\lambda)\varrho_1 + \lambda\varrho_2)$ is achievable under MI leakage, while the pair $(D_\lambda = (1-\lambda)D_1 + \lambda D_2, \varrho_\lambda = \log_2[(1-\lambda)2^{\varrho_1} + \lambda 2^{\varrho_2}])$ is achievable for MaxL.*

B. Connection to Rate-Distortion Theory

Consider an information source sequence with independent and identically distributed 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.

To relate the leakage-download function to the rate-distortion function, the desired file index m needs to be added as an argument to the answer-length function L , without changing the download cost in (4). Then, the leakage and the download cost play similar roles as the compression rate and the average distortion, respectively. This can be achieved by letting $L(m, \mathbf{q}) \triangleq L(\mathbf{q})$ 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). Below, we will equivalently use either $L(\mathbf{q})$ or $L(m, \mathbf{q})$ (as defined above).

IV. PARTITION WPIR SCHEME

In [11], a WPIR scheme based on partitioning was proposed. The scheme is formally described as follows. The files are first partitioned into η equally-sized partitions, each consisting of M_η files, where $M_\eta = M/\eta \in \mathbb{N}$. Assume that the requested file $\mathbf{X}^{(m)}$ belongs to the j -th partition, where $j \in [\eta]$. Then, the query \mathbf{Q} is constructed as

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

where $\tilde{\mathbf{Q}}$ is the query of an existing M_η -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_η -file WPIR scheme \mathcal{C} with achievable download-leakage pair $(\tilde{D}, \tilde{\varrho})$ exists. Then, the download-leakage pair*

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

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

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 $\mathcal{C}^{\text{basic}}$ (it achieves the pair in (7) with $\tilde{\varrho} = 0$). 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 VI, where for any subset $\mathcal{M} \subseteq [M]$ that includes the requested file, all files in \mathcal{M} are downloaded.

V. 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

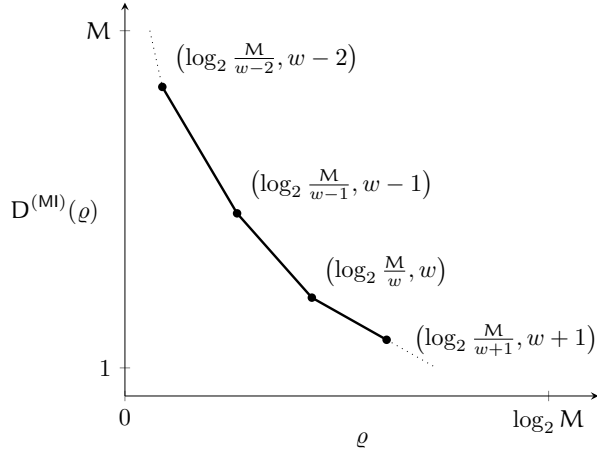


Fig. 1. An illustration of the function $D^{(MI)}(\varrho)$, which is defined by many linear functions.

of files for the MI and MaxL privacy metrics. The optimal download-leakage region for the MI privacy metric is stated in the following theorem. For the sake of illustration, we consider the normalized leakage-download function $\bar{\rho}^{(\cdot)} \triangleq \frac{\rho^{(\cdot)}}{\log_2 M}$.

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

$$C^{(MI)}(\bar{\varrho}) = \left[w + \frac{\log_2 \frac{M}{w}}{\log_2 \frac{w}{w-1}} - \frac{\bar{\varrho} \log_2 M}{\log_2 \frac{w}{w-1}} \right]^{-1},$$

for $1 - \frac{\log_2 w}{\log_2 M} \leq \bar{\varrho} \leq 1 - \frac{\log_2(w-1)}{\log_2 M}$, $w \in [2 : M]$.

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

$$C^{(MaxL)}(\bar{\varrho}) = \left[(2w-1) - \frac{2\bar{\varrho} \log_2 M}{\frac{M}{w-1} - \frac{M}{w}} \right]^{-1},$$

for $1 - \frac{\log_2 w}{\log_2 M} \leq \bar{\varrho} \leq 1 - \frac{\log_2(w-1)}{\log_2 M}$, $w \in [2 : M]$.

It is worthwhile noting that the download-leakage function $[C^{(\cdot)}(\cdot)]^{-1}$ is a piecewise continuous function. For the MI privacy metric, $[C^{(MI)}(\varrho)]^{-1}$ is a piecewise linear function as demonstrated in Fig. 1 (without normalization). Note also that, when $M_\eta = M/\eta \in \mathbb{N}$, the basic scheme $\mathcal{C}^{\text{basic}}$ 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 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 upperbounded as

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

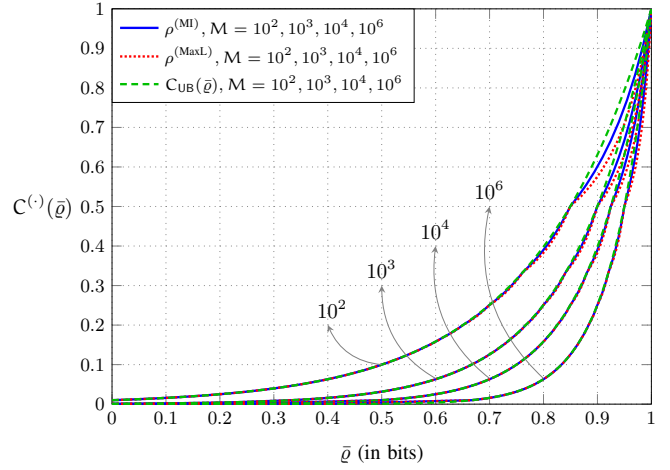


Fig. 2. The capacity $C^{(\cdot)}(\bar{\varrho})$ and its upper bound $C_{UB}(\bar{\varrho})$ for a number of files $M = 10^2, 10^3, 10^4, 10^6$ with privacy metrics $\rho^{(MI)}$ and $\rho^{(MaxL)}$.

In Fig. 2, the capacity $C^{(\cdot)}(\bar{\varrho})$ and the upper bound $C_{UB}(\bar{\varrho})$ are plotted for $M = 10^2, 10^3, 10^4, 10^6$, which illustrates the asymptotic behavior of $C^{(\cdot)}(\bar{\varrho})$ as M tends to infinity. Observe that from Theorems 2 and 3, 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

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

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

VI. ACHIEVABILITY OF THEOREMS 2 AND 3

Throughout this section, for simplicity, we set $\beta = 1$, i.e., $H(\mathbf{X}^{(m)}) = 1, \forall m \in [M]$. In fact, our proposed single-server WPIR capacity-achieving scheme works for an arbitrary file size β , which indicates that subpacketization does not improve the performance of single-server WPIR.

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 $M = 3$. We start by considering three achievable download-leakage pairs $(D_1, \varrho_1) = (1, \log_2 3)$, $(D_2, \varrho_2) = (2, \log_2 \frac{3}{2})$, and $(D_3, \varrho_3) = (3, \log_2 \frac{3}{3} = 0)$. In terms of the MI or MaxL privacy metrics, one can check that the pair (D_w, ϱ_w) is achieved by the conditional distribution $P_{Q_w|M}$, $w \in [3]$, listed in Table I. Moreover, it is clear that the retrievability condition in (1) is satisfied for any convex combination of the three PMFs.

Now, construct two conditional query distributions as follows,

$$P_{Q_{\lambda_1}|M} = (1 - \lambda_1)P_{Q_2|M} + \lambda_1 P_{Q_1|M}, \quad (8)$$

$$P_{Q_{\lambda_2}|M} = (1 - \lambda_2)P_{Q_3|M} + \lambda_2 P_{Q_2|M}, \quad (9)$$

TABLE I
THE CONDITIONAL PMFS $P_{\mathbf{Q}_w|M}$, $w \in [3]$.

\mathcal{Q}_1	$P_{\mathbf{Q}_1 M}(q 1)$	$P_{\mathbf{Q}_1 M}(q 2)$	$P_{\mathbf{Q}_1 M}(q 3)$	\mathbf{A}	$P_{\mathbf{Q}_1}(q)$
(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	$X_1^{(3)}$	$\frac{1}{3}$

\mathcal{Q}_2	$P_{\mathbf{Q}_2 M}(q 1)$	$P_{\mathbf{Q}_2 M}(q 2)$	$P_{\mathbf{Q}_2 M}(q 3)$	\mathbf{A}	$P_{\mathbf{Q}_2}(q)$
(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	$\frac{1}{2}$	$\frac{1}{2}$	$\{X_1^{(2)}, X_1^{(3)}\}$	$\frac{1}{3}$

\mathcal{Q}_3	$P_{\mathbf{Q}_3 M}(q 1)$	$P_{\mathbf{Q}_3 M}(q 2)$	$P_{\mathbf{Q}_3 M}(q 3)$	\mathbf{A}	$P_{\mathbf{Q}_3}(q)$
(1, 1, 1)	1	1	1	$\{X_1^{(1)}, X_1^{(2)}, X_1^{(3)}\}$	1

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

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 (8)–(9) 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

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

and the MaxL

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

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 $M = 3$.

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

$$(D_w, \varrho_w) = \left(w, \log_2 \frac{M}{w} \right), \quad w \in [M],$$

are achievable.

1) *Query Generation*: Consider M random queries \mathbf{Q}_w , $w \in [M]$, whose alphabet is $\mathcal{Q}_w \triangleq \{\mathbf{q} = (q_1, \dots, q_M) \in \{0, 1\}^M : w_H(\mathbf{q}) = w\}$. Each query $\mathbf{q} \in \mathcal{Q}_w$ sent to the server is generated by the conditional PMF

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

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

2) *Answer Construction*: The answer function φ 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 $\mathbb{E}_{P_M P_{\mathbf{Q}_w|M}}[L(\mathbf{Q}_w)] = w$. The MI leakage is

$$\rho^{(\text{MI})}(P_{\mathbf{Q}_w|M}) = I(M; \mathbf{Q}_w) = \log_2 M - \log_2 w = \log_2 \frac{M}{w}$$

and the MaxL is

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

Notice that the presented capacity-achieving WPIR scheme can be seen as a generalization of the basic WPIR scheme $\mathcal{C}^{\text{basic}}$. If $M/\eta = w \in \mathbb{N}$, it can be shown that the download-leakage pair $(w, \log_2 \eta) = (w, \log_2 \frac{M}{w})$ is also achievable by $\mathcal{C}^{\text{basic}}$ for both the MI and MaxL privacy metrics.

VII. CONVERSE OF THEOREMS 2 AND 3

The converse proofs for the MI and MaxL privacy metrics are fully elaborated in the extended version [16, Secs. VII and VIII]. Here, we briefly outline the main steps of the converse proof for MI leakage.

Consider a single-server WPIR scheme where the leakage at the server is measured by $I(M; \mathbf{Q})$. A general converse for Theorem 2 can be derived from the download-leakage function of a given leakage constraint ϱ , or equivalently, from the leakage-download function of a given download cost constraint D . Similar to (5), the leakage-download function can be formulated by the convex minimization problem

$$P_{\mathbf{Q}|M} : \min_{\mathbb{E}[L(M, \mathbf{Q})] \leq D} I(M; \mathbf{Q}). \quad (10)$$

The proof consists of two parts as described below.

Part 1: With some technical efforts, we can prove that (10) is bounded from below by

$$\rho^{(\text{MI})}(D) \triangleq \min_{P_{\mathbf{U}|M} : \mathbb{E}[L(M, \mathbf{U})] \leq D} I(M; \mathbf{U}), \quad (11)$$

where \mathbf{U} is a length- M binary random vector and the length function $L(m, \mathbf{u})$ is constructed for any $\mathbf{u} \in \{0, 1\}^M$ by

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

Part 2: Since (11) can be related to the rate-distortion function with a certain distortion measure, we use a useful result from the rate-distortion theory (see [16, Lem. 3]) to prove a lower bound on (11). It can be shown that the pair $(\bar{\varrho}, D^{(\text{MI})}(\bar{\varrho}))$ of Theorem 2 lies on the lower bound of (11).

VIII. CONCLUSION

We characterized the capacity of single-server WPIR with an arbitrary number of files for the MI and MaxL privacy metrics. We showed that, interestingly, the optimal download rate is higher than the single-server PIR capacity. Moreover, we showed that if the server can not determine the identity of the requested file, then the capacity tends to zero as the number of files goes to infinity.

REFERENCES

- [1] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan, "Private information retrieval," in *Proc. 36th Annu. IEEE Symp. Found. Comp. Sci. (FOCS)*, Milwaukee, WI, USA, Oct. 23–25, 1995, pp. 41–50.
- [2] 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.
- [3] 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.
- [4] 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.
- [5] 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. 25–29, 2018, pp. 1–5.
- [6] 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.
- [7] 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.
- [8] 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.
- [9] L. Holzbaur, R. Freij-Hollanti, J. Li, and C. Hollanti, "Towards the capacity of private information retrieval from coded and colluding servers," Mar. 2019, arXiv:1903.12552v5 [cs.IT].
- [10] R. R. Toledo, G. Danezis, and I. Goldberg, "Lower-cost ϵ -private information retrieval," in *Proc. Privacy Enhancing Technol. Symp. (PETS)*, Darmstadt, Germany, Jul. 19–22, 2016, pp. 184–201.
- [11] 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. 7–12, 2019, pp. 1257–1261.
- [12] 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. 7–12, 2019, pp. 1262–1266.
- [13] G. Smith, "On the foundations of quantitative information flow," in *Proc. 12th Int. Conf. Found. Softw. Sci. Comput. Struct. (FoSSaCS)*, York, U.K., Mar. 22–29, 2009, pp. 288–302.
- [14] G. Barthe and B. Köpf, "Information-theoretic bounds for differentially private mechanisms," in *Proc. 24th Comput. Secur. Found. Symp. (CSF)*, Cernay-la-Ville, France, Jun. 27–29, 2011, pp. 191–204.
- [15] 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.
- [16] H.-Y. Lin, S. Kumar, E. Rosnes, A. Graell i Amat, and E. Yaakobi, "The capacity of single-server weakly-private information retrieval," Jan. 2020, arXiv:2001.08727v1 [cs.IT].